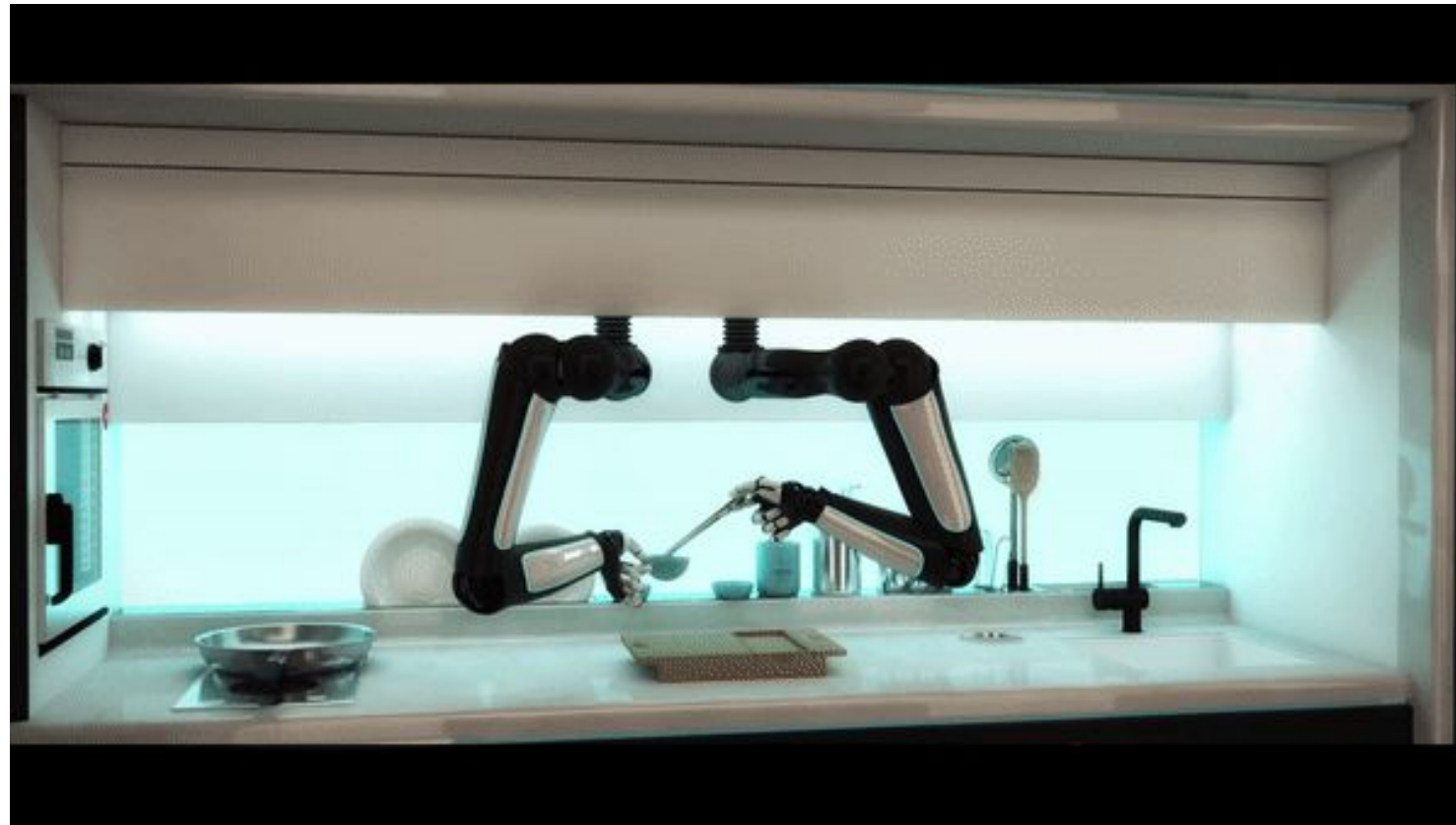


Guiding Policies with Language via Meta-Learning

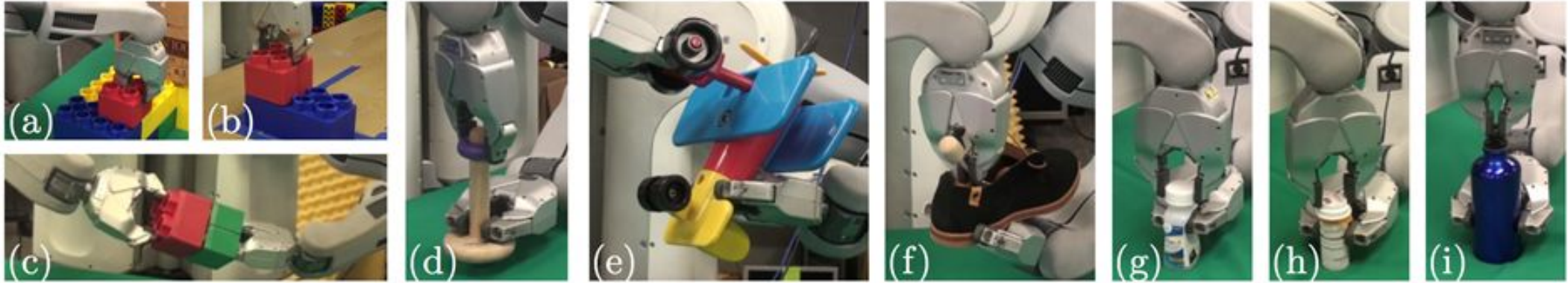
John D (JD) Co-Reyes, Abhishek Gupta, Suvansh Sanjeev, Nick Altieri,
John DeNero, Pieter Abbeel, Sergey Levine



Ideal Robot



Learning new tasks quickly

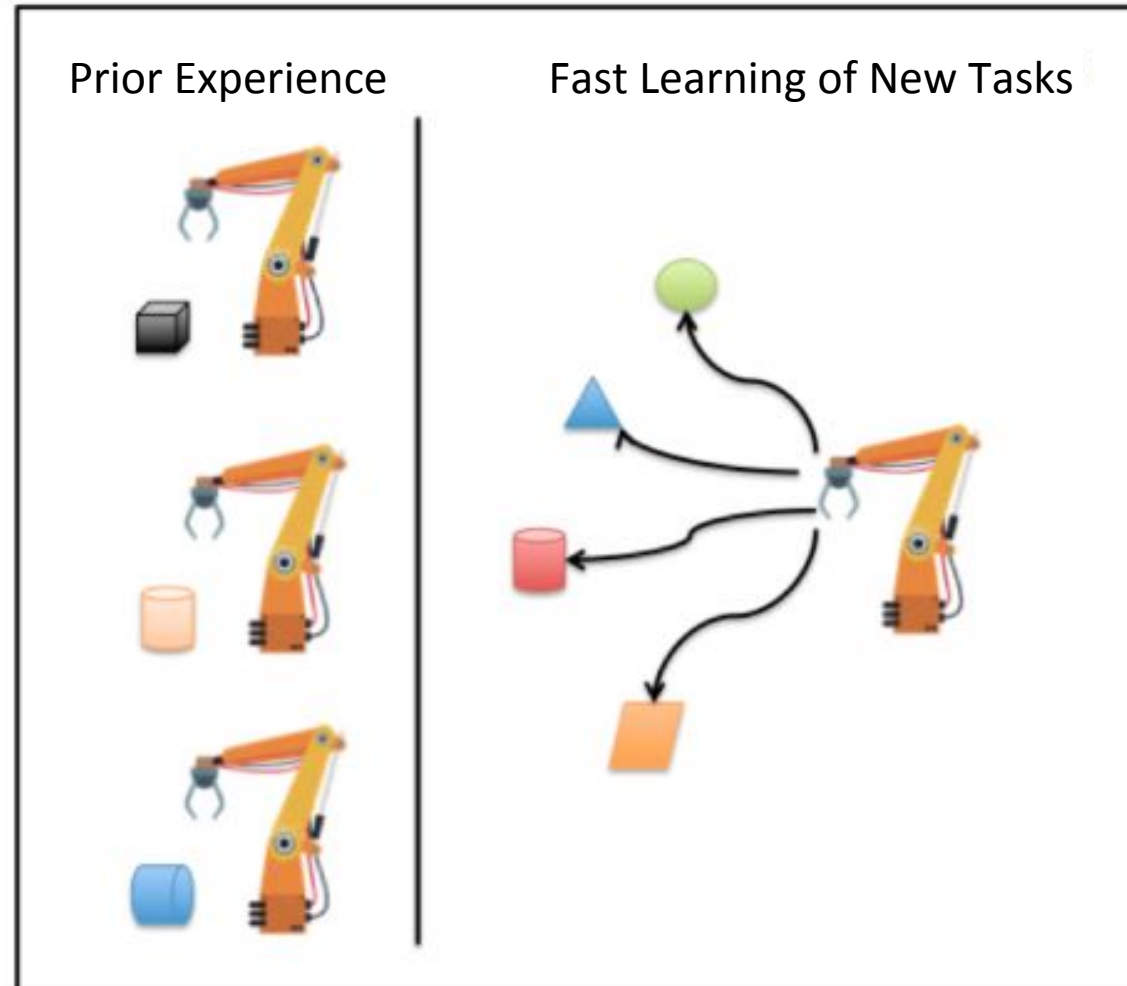


- Want diverse range of skills
- Cost of supervision can be high
- Want to learn new things with as little supervision as possible

Meta-RL

Leverage prior experience to quickly learn new tasks

Meta-Training



Meta-Testing

Problem with reward design

Hard to design



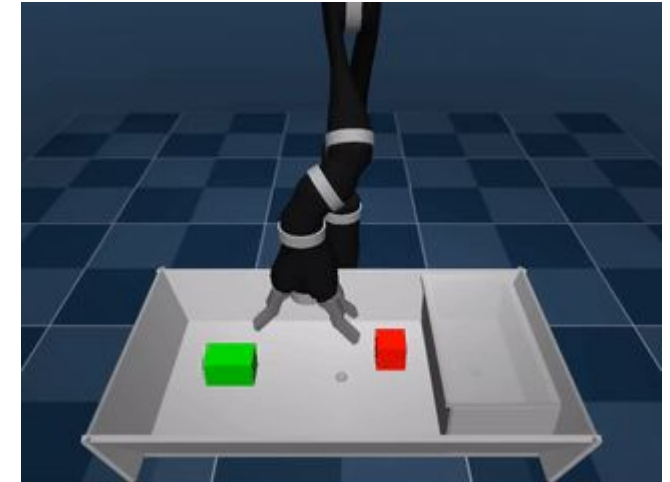
$r = \text{game score}$

Hard to provide



$$\begin{aligned} r_{shaped}(s, a) = & \|hammer - palm\|_2 \\ & + \|hammer - nail\|_2 \\ & + \|nail - goal\| \end{aligned}$$

Hard to learn from



$r = \mathbb{1}(\text{both blocks in box})$

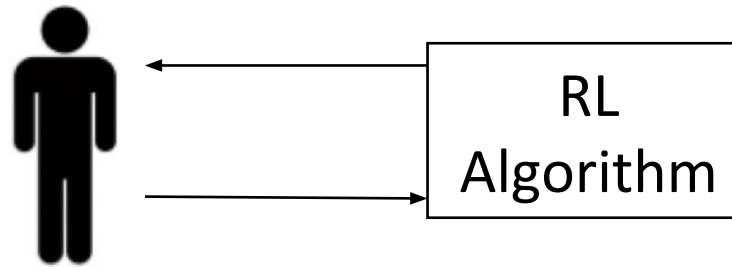
More natural way to provide supervision

Human feedback

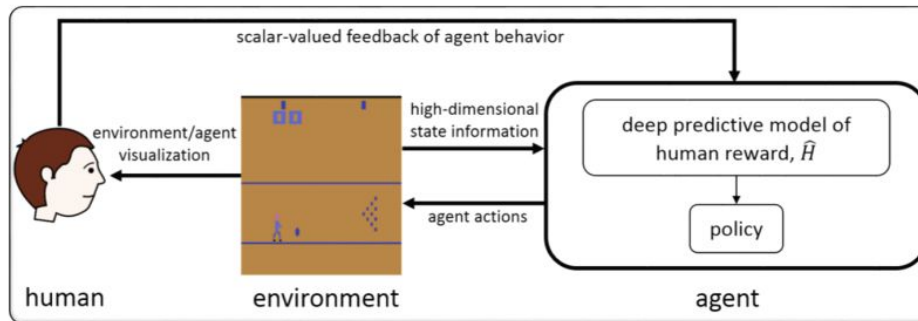


Human-in-the-loop supervision

Replace reward with human feedback

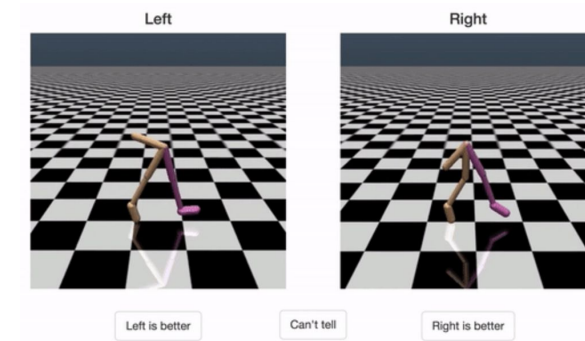


Deep TAMER



Warnell et al

Preferences

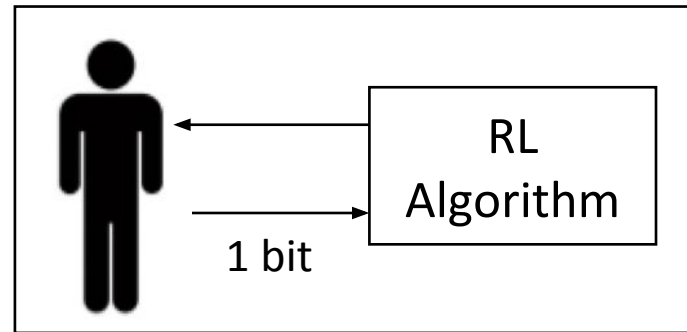


Christiano et al

Why current methods are insufficient?

Very few bits of information per intervention

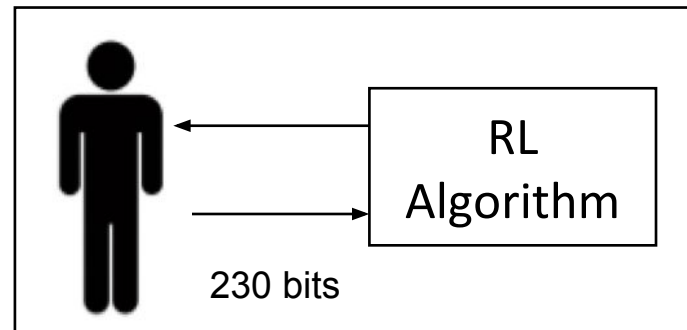
Scalar Feedback



Significant
human effort

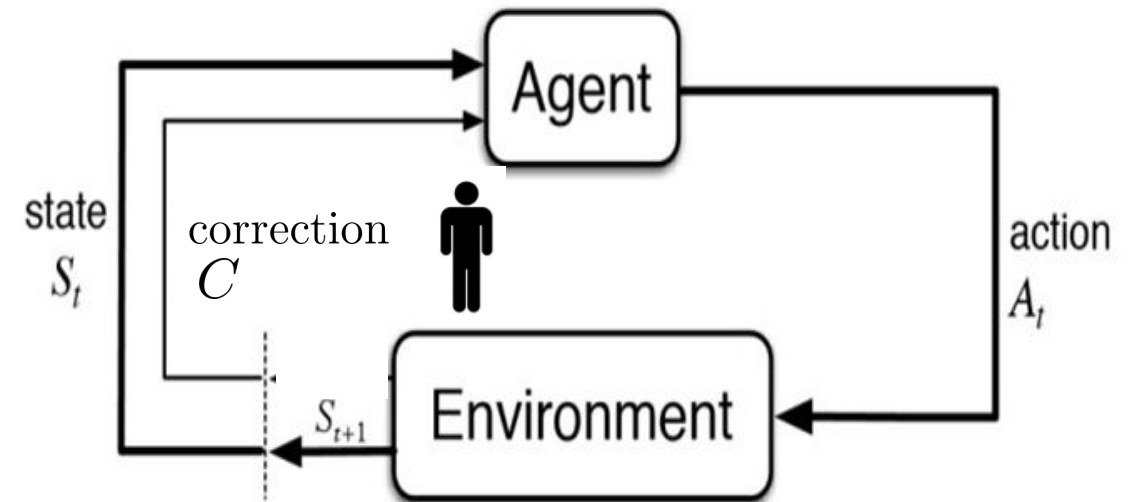
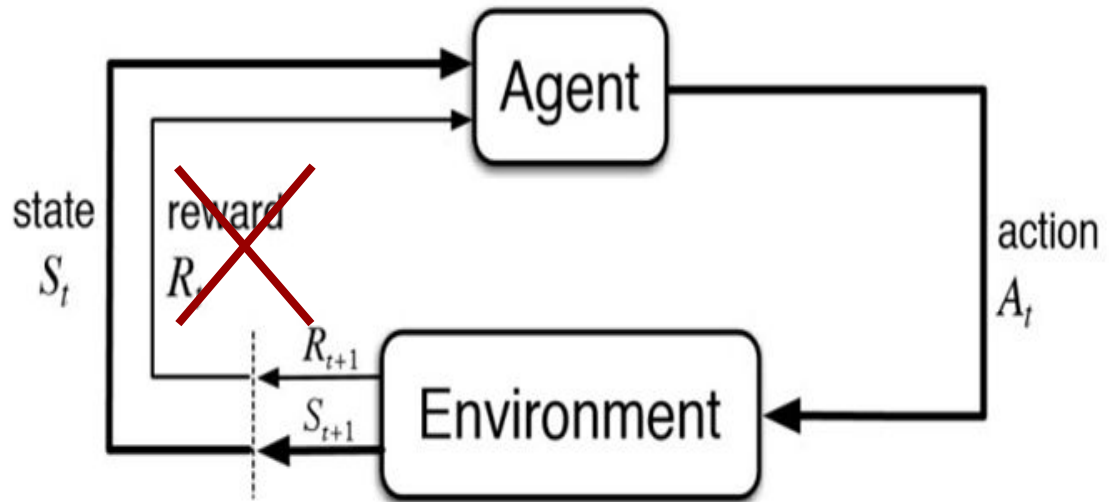
More bits of information per intervention

Language Feedback



Less
human effort

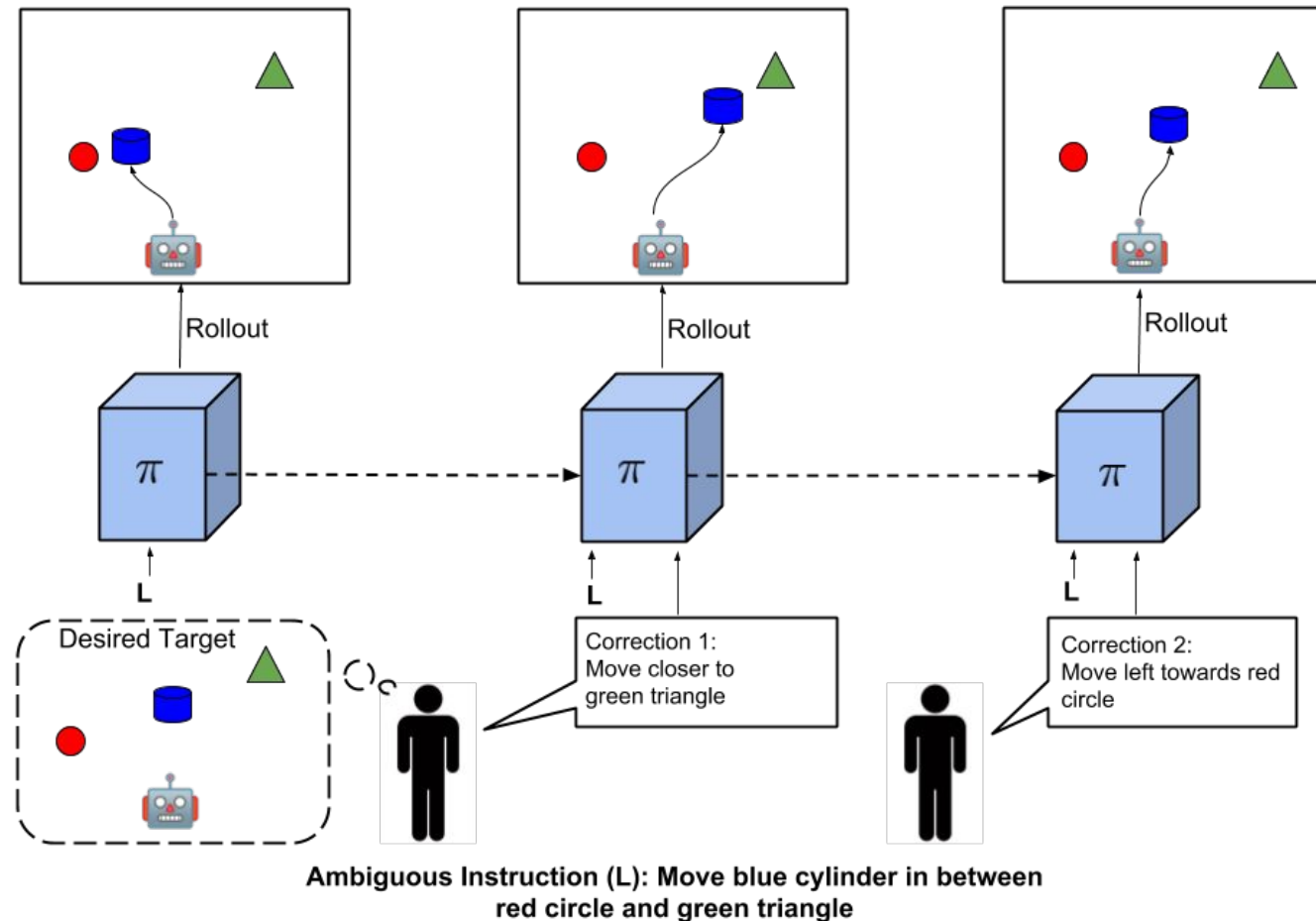
Language Corrections



Problem Setting

Agent provided with ambiguous/incomplete instruction

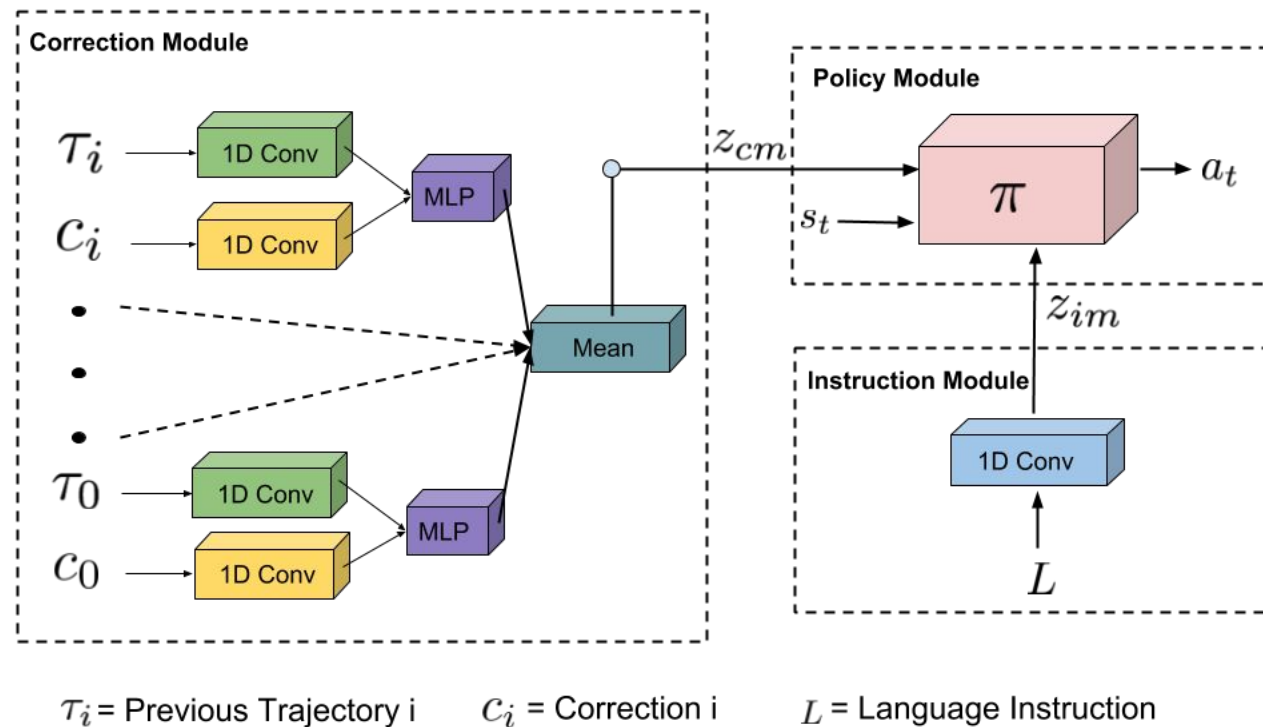
Quickly incorporate language corrections in the loop



Language Guided Policy Model

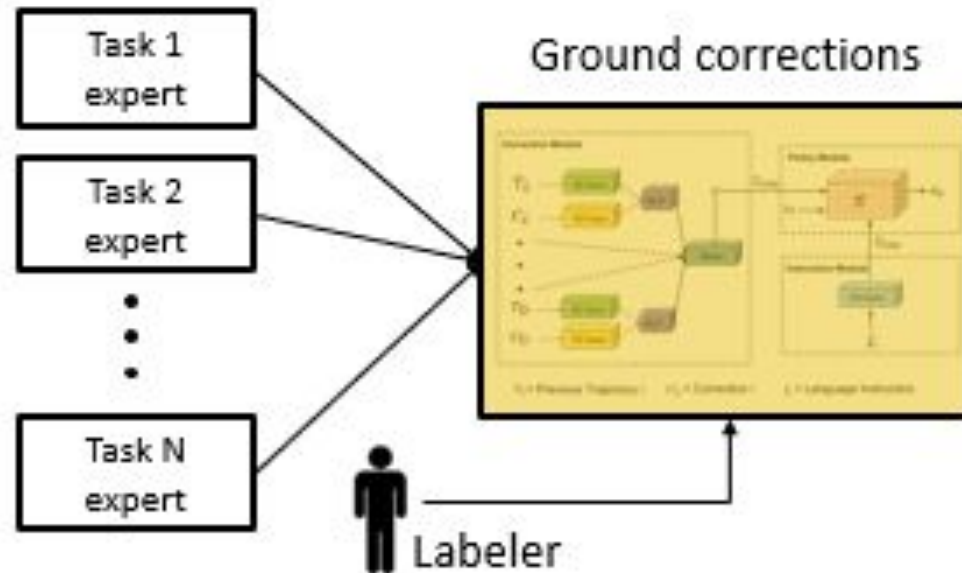
Model improves based on previous trajectories and corrections.

3 modules – corrections, policy and instruction modules

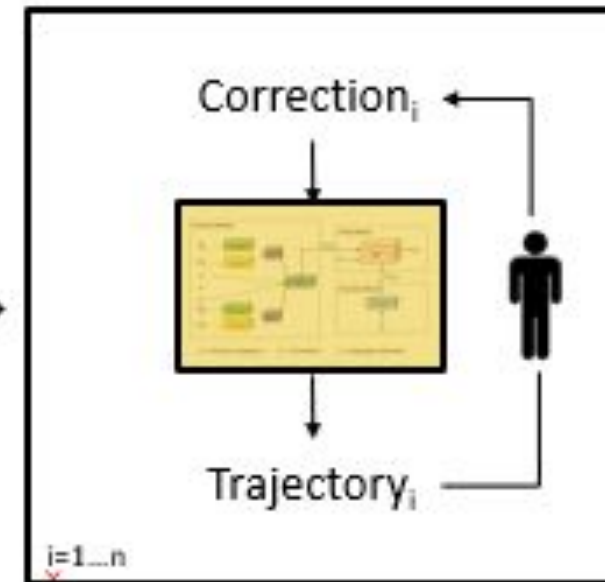


Algorithm Overview

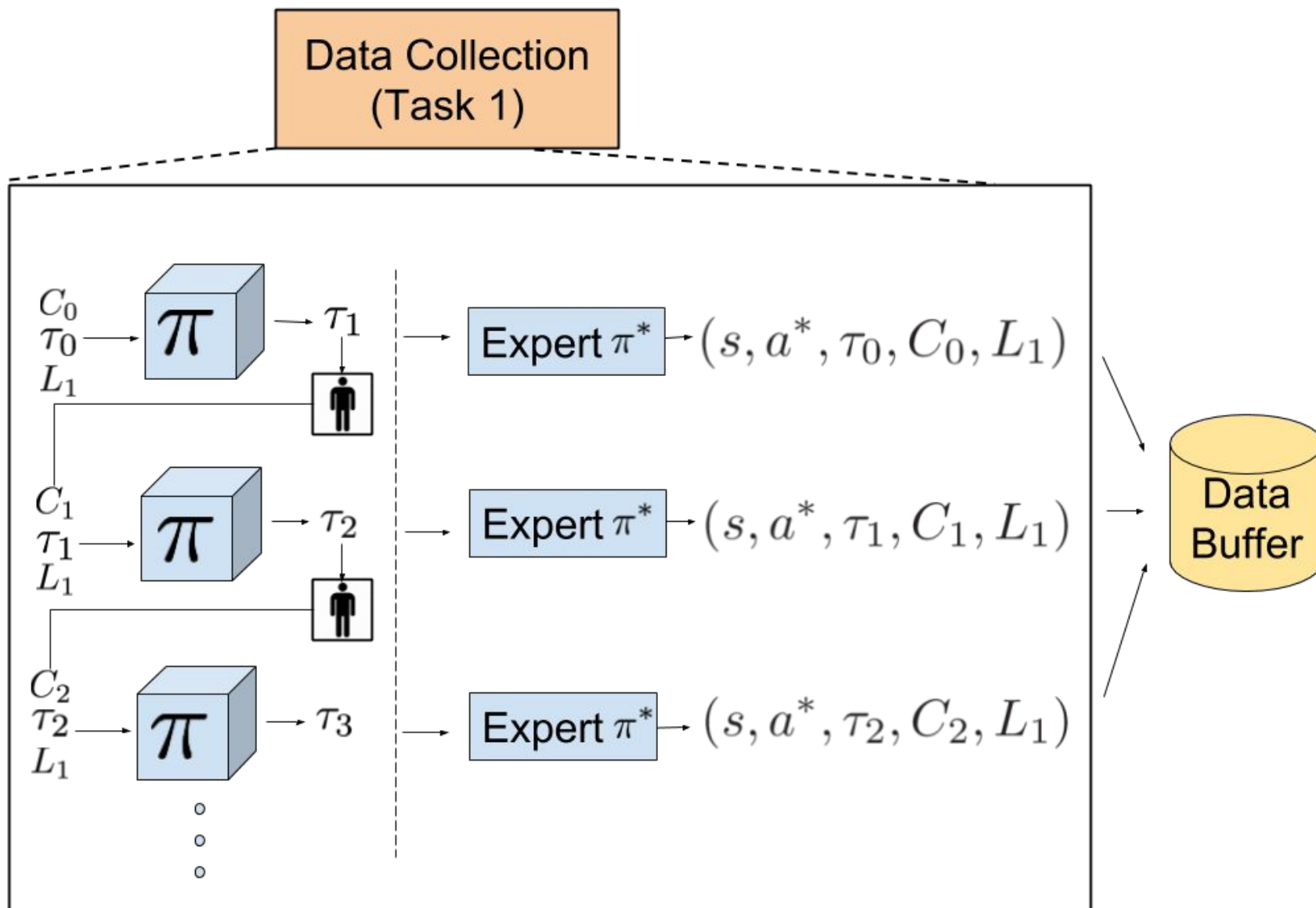
Meta-Training



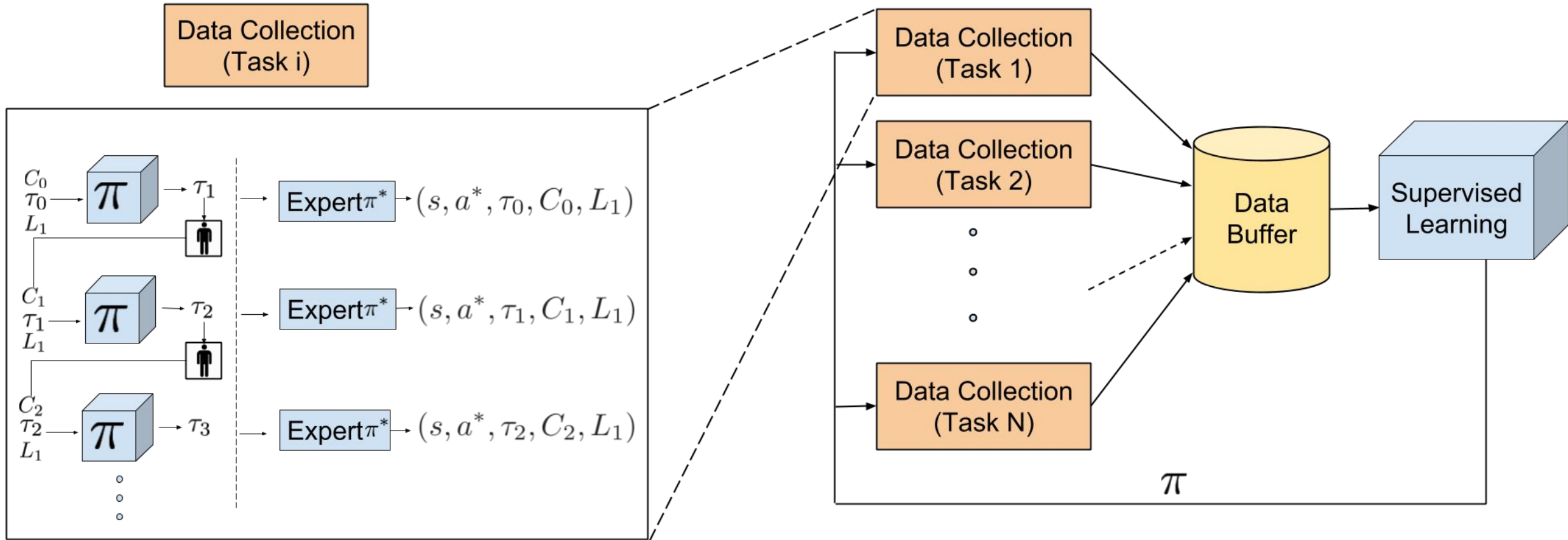
Meta-Testing



Meta-Training

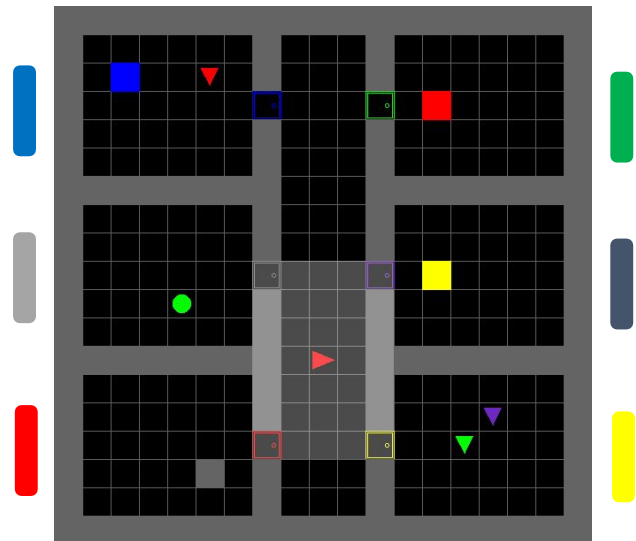
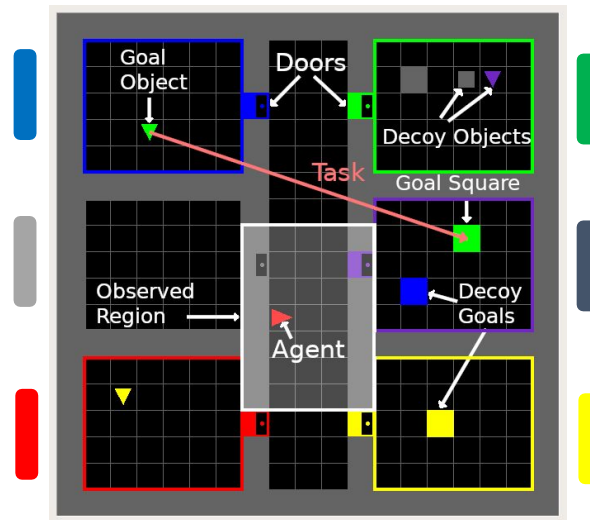


Meta-Training

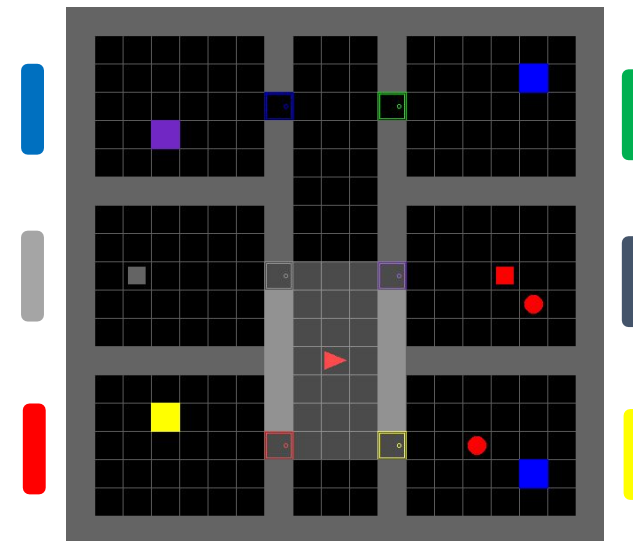


Experimental Setup

Multi-room domain



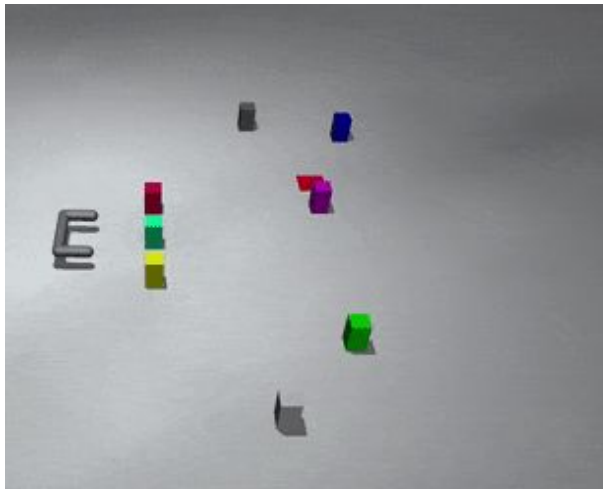
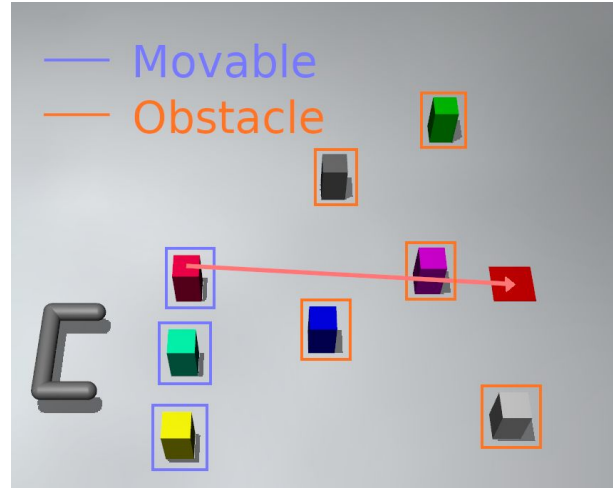
Instruction: Move green triangle to yellow goal.



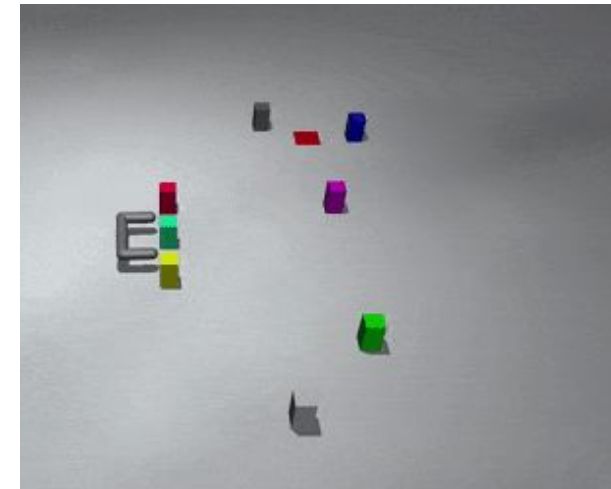
Instruction: Move red square to yellow goal.

Experimental Setup

Block pushing domain



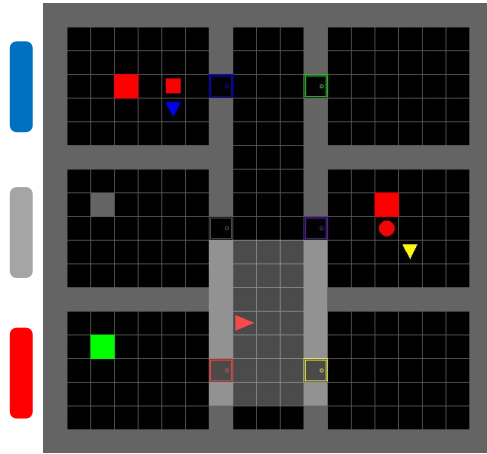
Instruction: Move red block above magenta block.



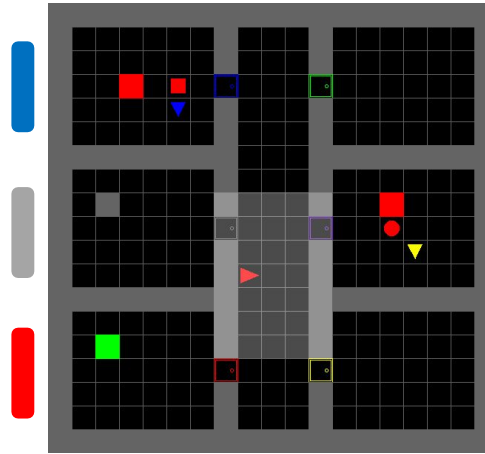
Instruction: Move cyan block left of blue block.

Quick Learning of New Tasks

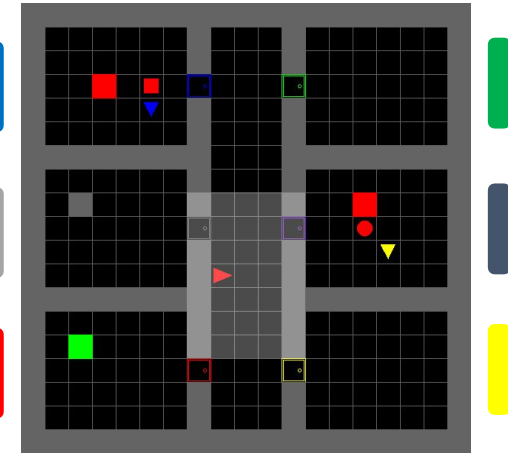
Instruction: Move blue triangle to green goal.



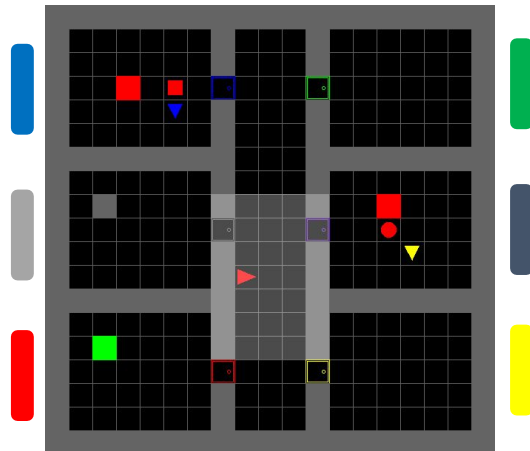
Correction 1: Enter the blue room.



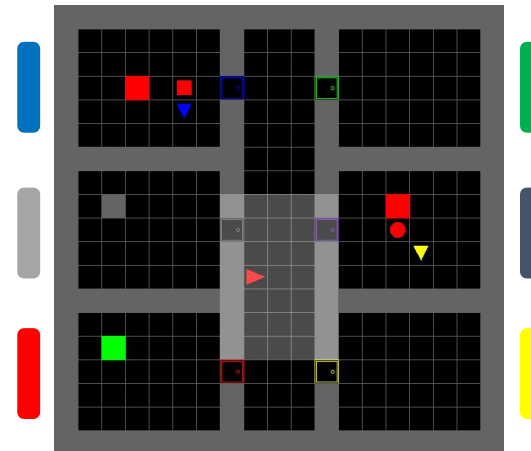
Correction 2: Enter the red room.



Correction 3: Exit the blue room.



Correction 4: Pick up the blue triangle.



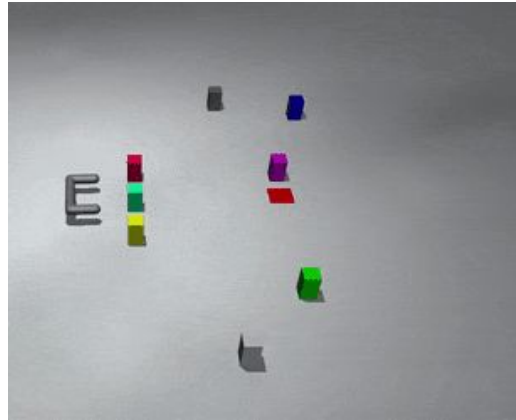
Solved

Quick Learning of New Tasks

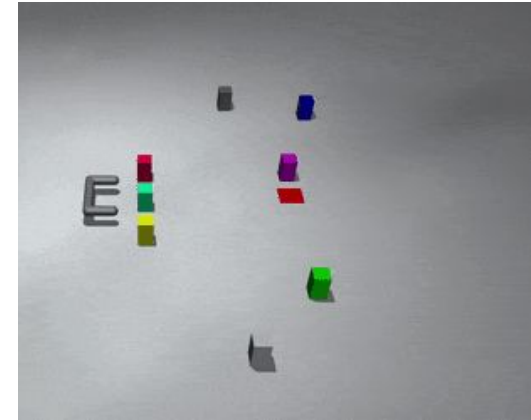
Instruction: Move cyan block below magenta block.



Correction 1: Touch cyan block.



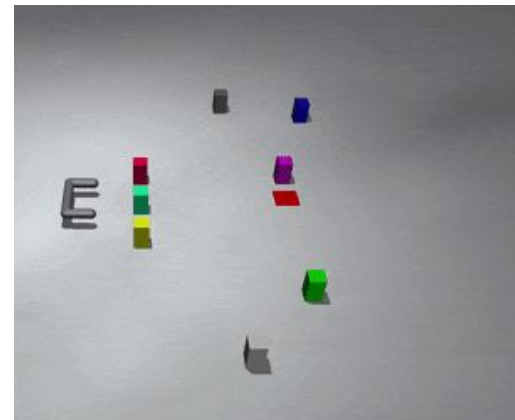
Correction 2: Move closer to magenta block.



Correction 3: Move a lot up.



Correction 4: Move a little up.



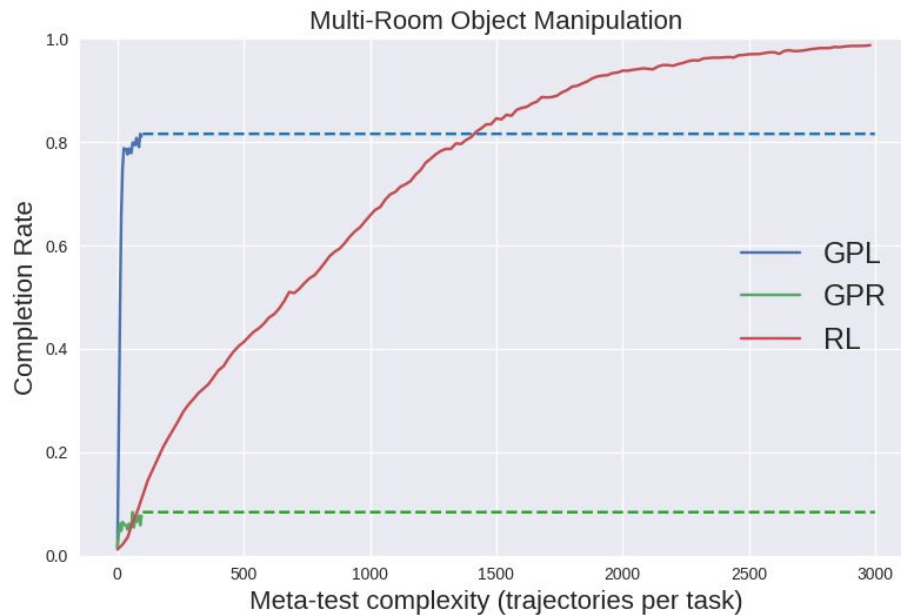
Solved

Quantitative Evaluation

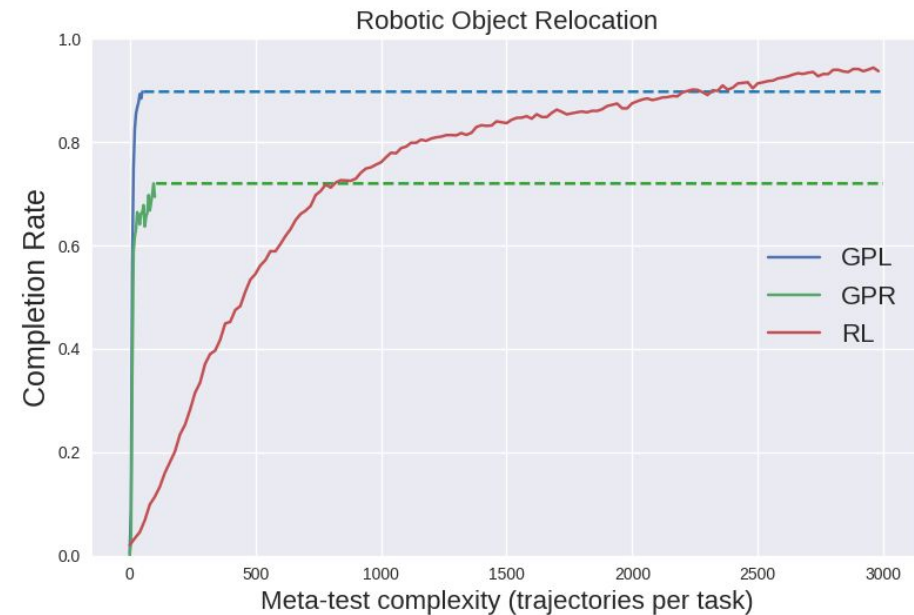
Success Rates on New Tasks

Env	Instruction	Full Info	MIVOA (Instr.)	MIVOA (Full Info)	c_0	c_1	c_2	c_3	c_4	c_5
Multi-room	0.075	0.73	0.067	0.63	0.066	0.46	0.65	0.73	0.77	0.82
Obj Relocation	0.64	0.96	0.65	-	0.65	0.80	0.84	0.85	0.88	0.90

Much quicker learning than using rewards



RL – PPO with reward used to train expert



GPL – Ours

Summary

- Avoid demos/reward functions using human-in-the-loop
- Language provides more information per intervention
- Ground language in multi-task setup; learn new tasks quickly with corrections

Thank you



Abhishek Gupta



Suvansh Sanjeev



Nick Altieri



John DeNero



Pieter Abbeel



Sergey Levine

