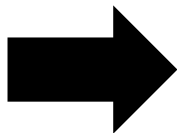# OpenAI

# Asymmetric self-play for automatic goal discovery in robotic manipulation
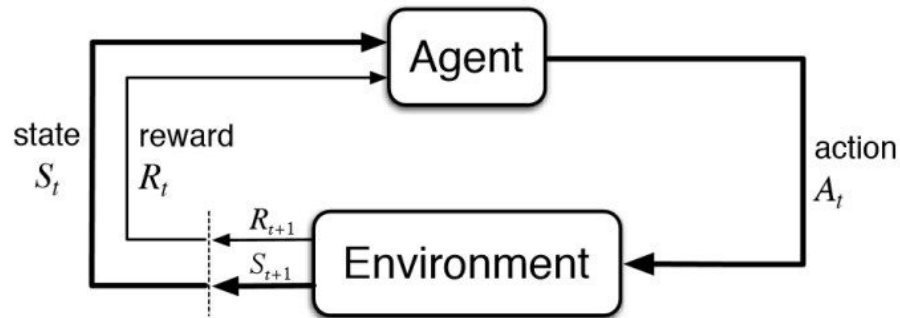
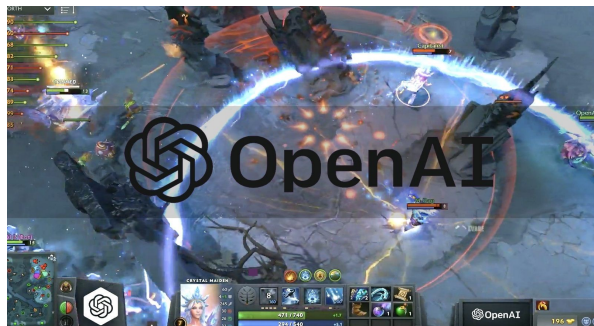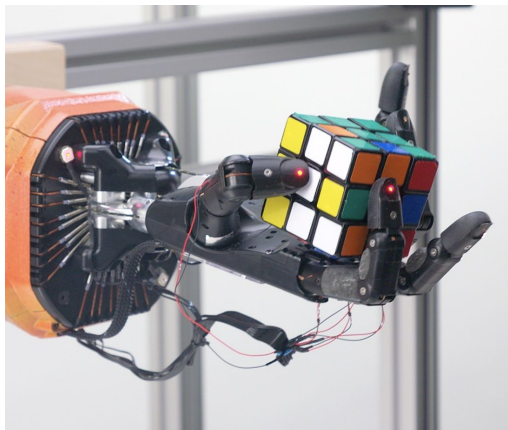lilian@openai, Apr 9 2021

**General Purpose Robot**

# Reinforcement Learning Basics



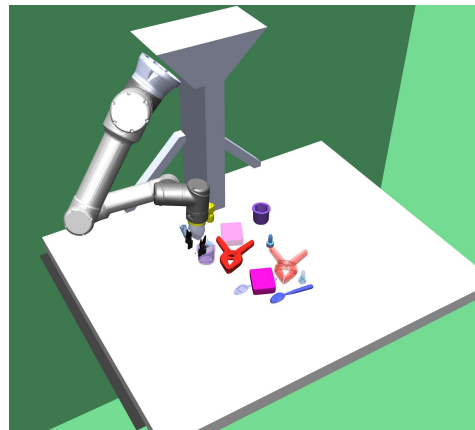Reinforcement Learning is **powerful**, but training needs **a lot of data**.

# Robotic Manipulation Tasks





**Solving rubik's cube with robot hand**

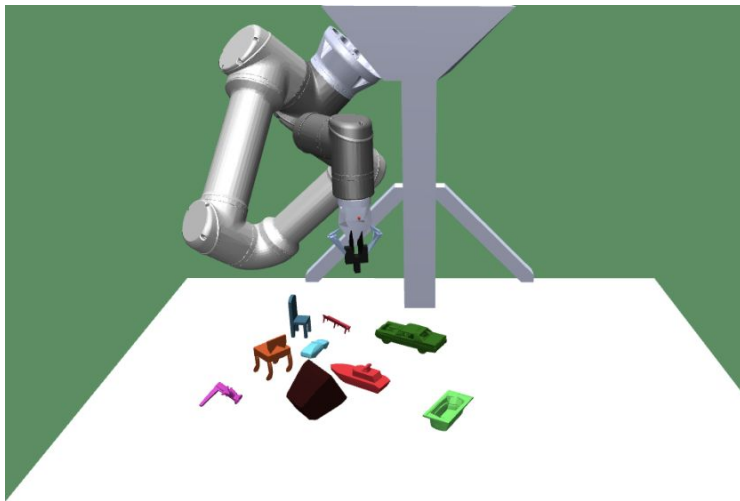The same RL control policy trained only in simulation can work in the real physical robot.

**Object rearrangement on the tabletop**

A single goal-conditioned policy can solve many manipulation tasks involving unseen arrangement and unseen objects.
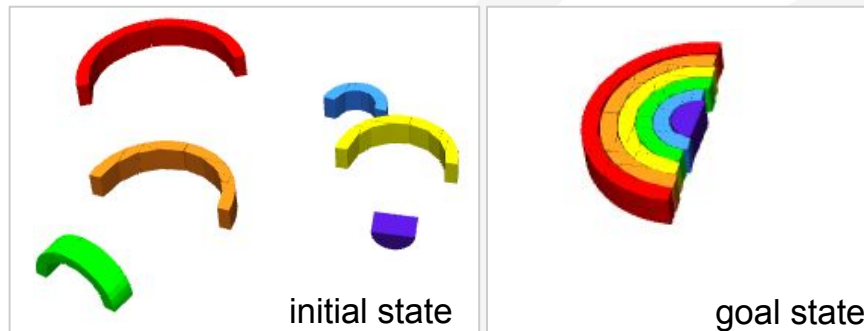
# Motivation

- Training a **single** goal-conditioned policy
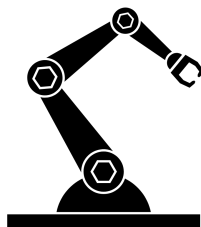- Solving **any** robotic manipulation task in an environment



Robotic manipulation environment:
one UR robot + gripper + table surface
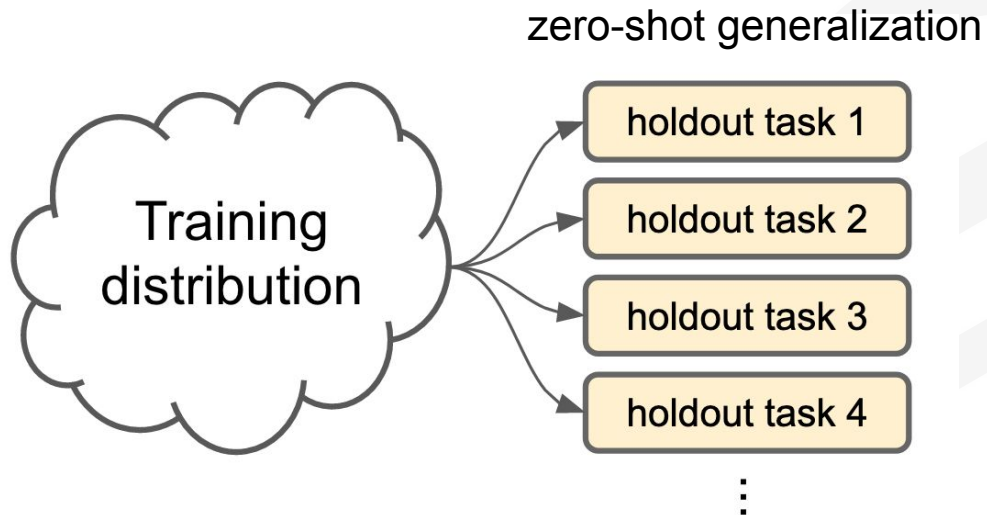


initial state
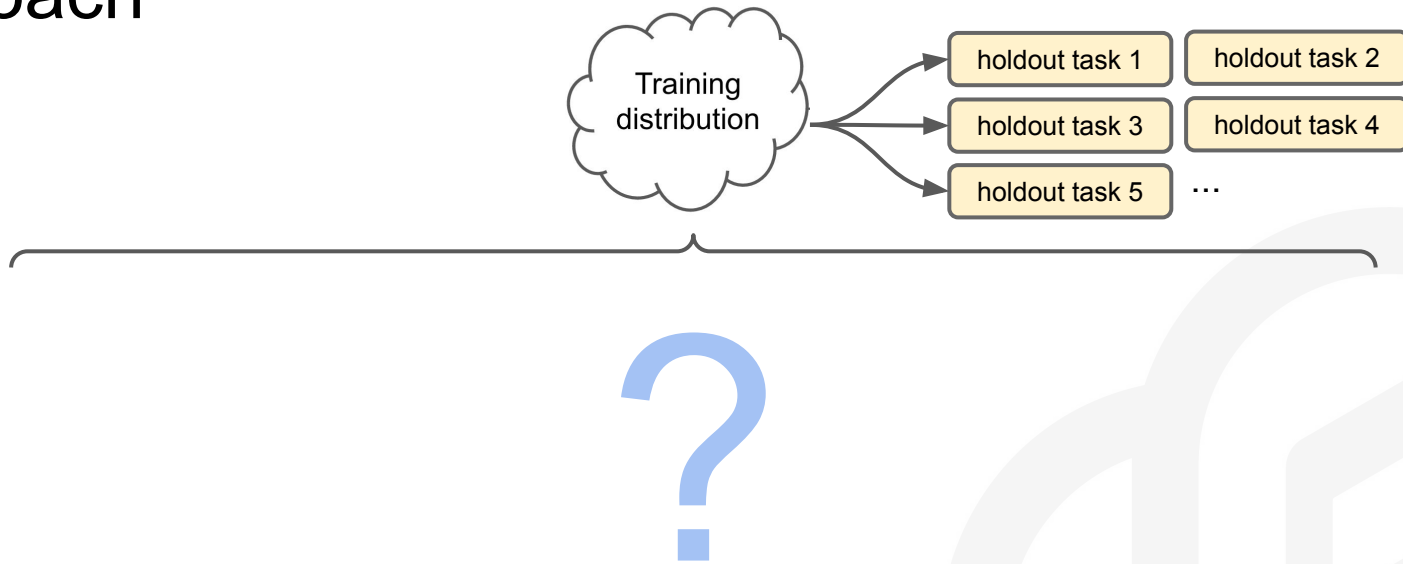
goal state

Task: Initial state → Goal state

# Approach

- Goal: One policy for all tasks
  - Training on a large training distribution (initial + goal states)
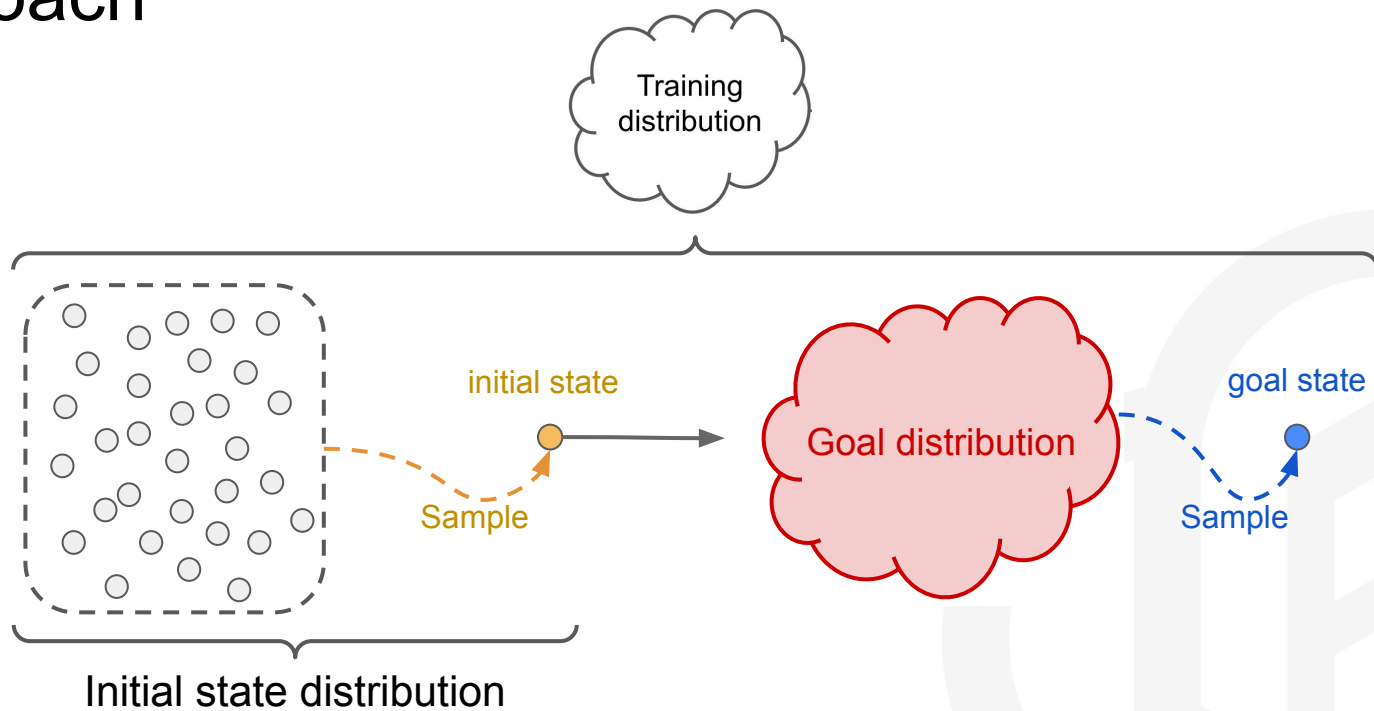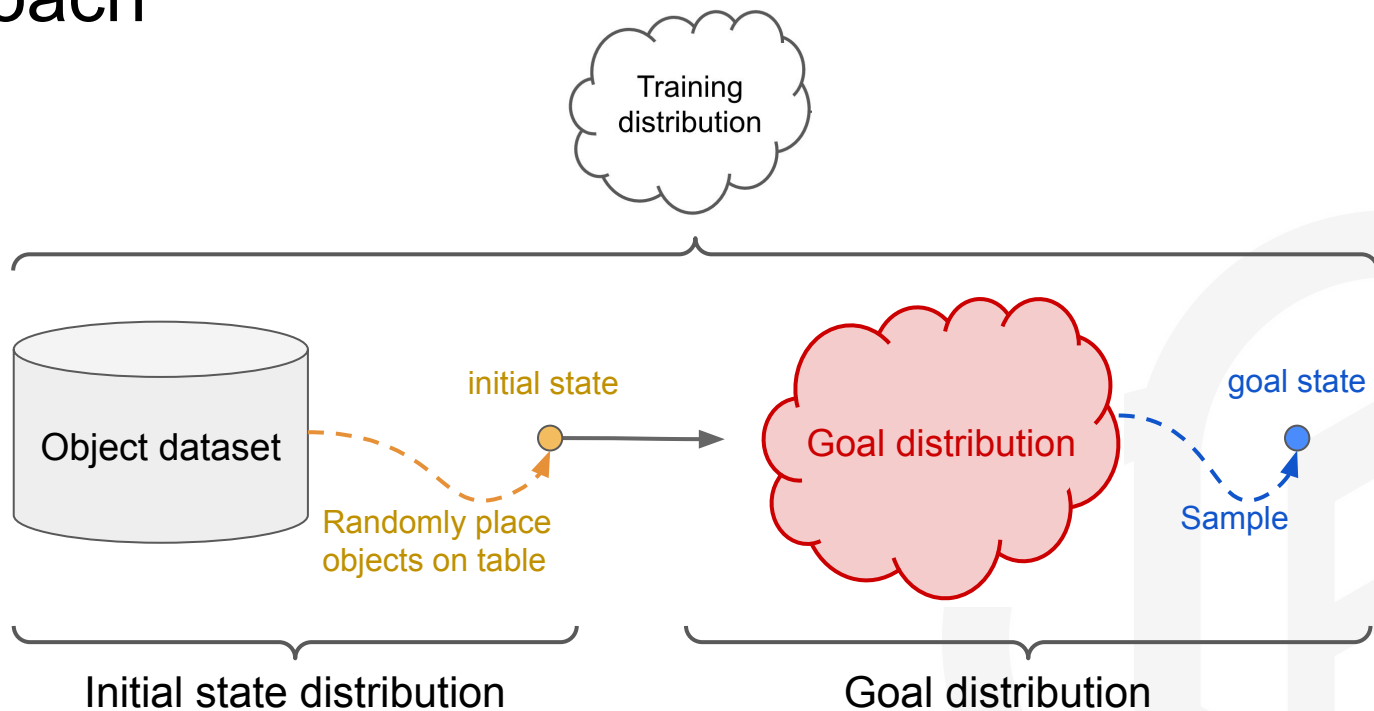  - Testing on unseen holdout tasks

zero-shot generalization



Training distribution

One policy

holdout task 1

holdout task 2

holdout task 3

holdout task 4

⋮

# Approach

# Approach

# Approach

# Approach



zero-shot generalization

holdout task 1    holdout task 2

holdout task 3    holdout task 4

holdout task 5    ...

Training distribution

Object dataset

initial state

Randomly place objects on table

Asymmetric self-play
[Sukhbaatar et al., 2018]

goal state

Sample

Initial state distribution

Goal distribution
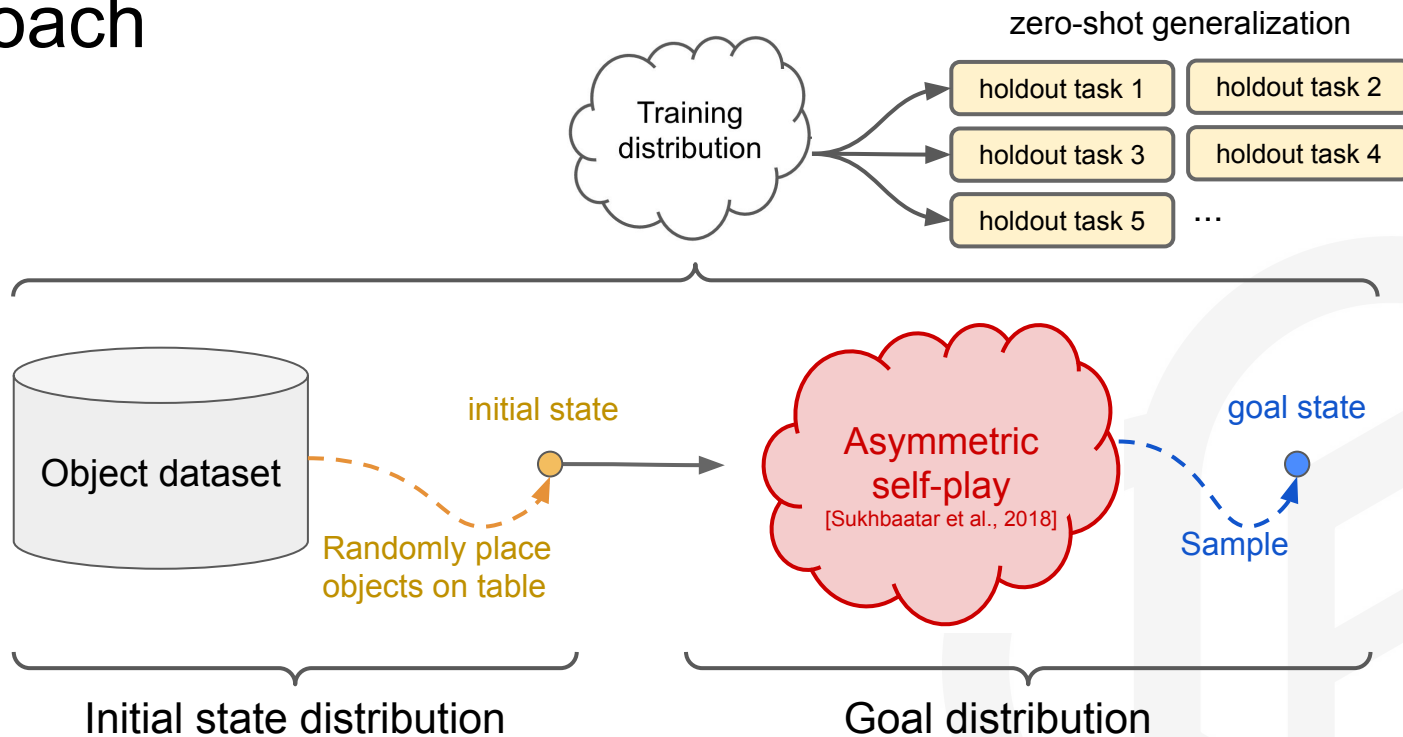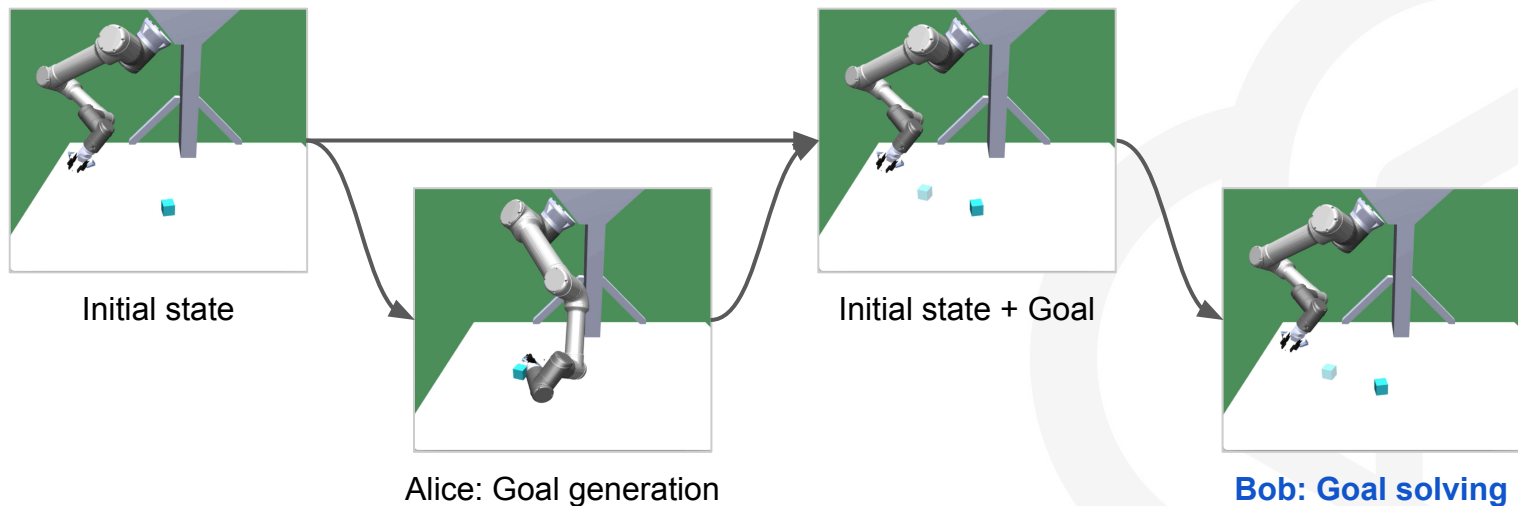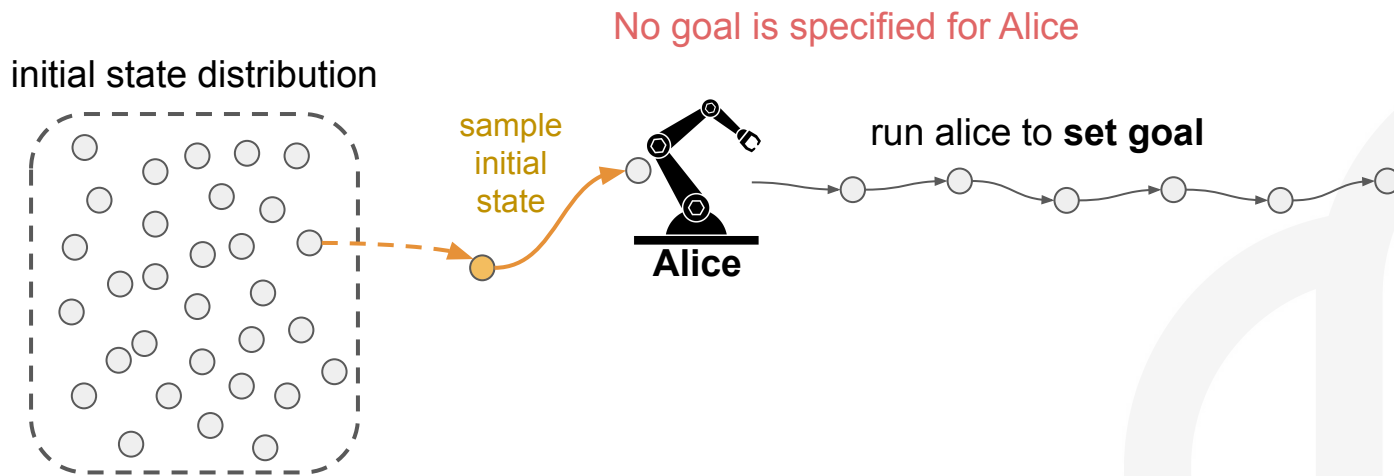
# Asymmetric Self-play for Robotics Manipulation

- Learning to generate goals + Learning to solve them:
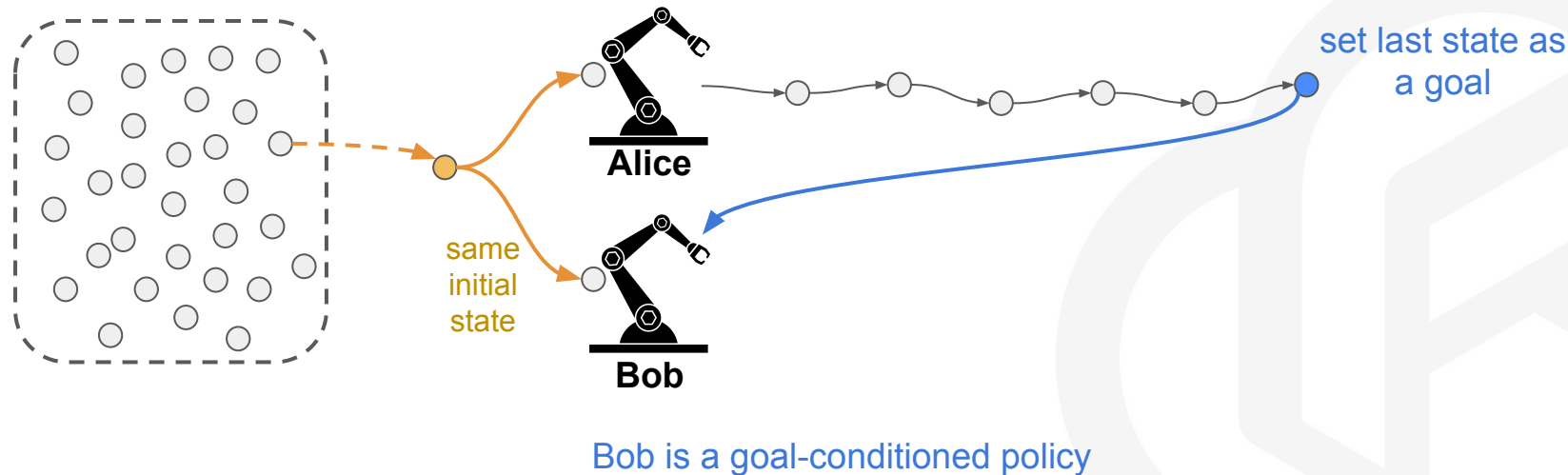  - Train two policies (Alice, Bob) for the same robotic hardware



Initial state

Alice: Goal generation

Initial state + Goal

**Bob: Goal solving**

# Asymmetric Self-play for Robotics Manipulation



No goal is specified for Alice

initial state distribution

sample initial state

run alice to **set goal**

**Alice**

# Asymmetric Self-play for Robotics Manipulation



initial state distribution

set last state as a goal

same initial state

**Alice**

**Bob**

Bob is a goal-conditioned policy

# Asymmetric Self-play for Robotics Manipulation

initial state distribution



set last state as a goal

Alice

same initial state

run bob to **solve goal**

Bob

Bob is a goal-conditioned policy

# Asymmetric Self-play for Robotics Manipulation



Incentivized to generate challenging goals

**Alice**

**Bob**

for Alice:
Positive reward
if Bob fails

for Bob: goal-conditioned reward

# Alice Behavioral Cloning (ABC)



$$\tau_A = \{(s_0, a_0, r_0), ..., (s_T, a_T, r_T)\}$$
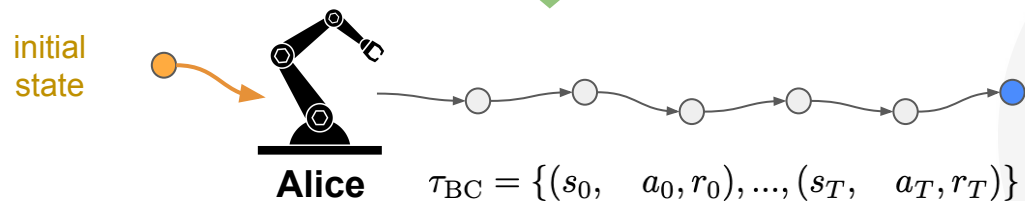
initial state

**Alice**

$s_T$

goal state

**Demonstration trajectory filtering**

Only if Bob fails this goal

(1) Copy this trajectory

initial state

**Alice**

$$\tau_{\mathrm{BC}} = \{(s_0, \quad a_0, r_0), ..., (s_T, \quad a_T, r_T)\}$$

# Alice Behavioral Cloning (ABC)



$$\tau_A = \{(s_0, a_0, r_0), ..., (s_T, a_T, r_T)\}$$

initial state

**Alice**

$s_T$

goal state

initial state

goal state $s_T$

**Alice**

(2) Relabeled goal solving trajectory

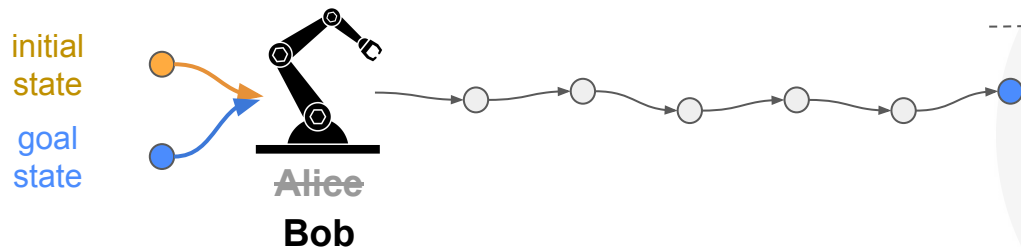$$\tau_{\mathrm{BC}} = \{(s_0, g, a_0, r_0), ..., (s_T, g, a_T, r_T)\}$$

$s_T$    $s_T$

# Alice Behavioral Cloning (ABC)

For Bob:

$$\mathcal{L} = \mathcal{L}_{\mathrm{RL}} + \beta \underbrace{\mathcal{L}_{\mathrm{abc}}}$$

Behavioral cloning loss

(3) Use demonstration

initial state

Alice

goal state

initial state

goal state

~~Alice~~
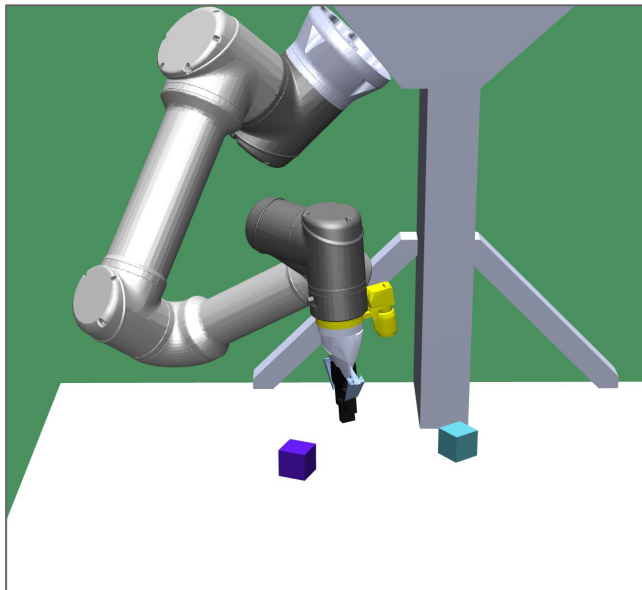
Bob

# Stabilizing Alice Behavioral Cloning (ABC)

- Demonstration filtering:
  - Collect demonstration only for failed goal
- PPO (Schulman et al., 2017)-style clipping:
  - Prevent drastic policy change

$$\mathcal{L}_{\mathrm{bc}} = -\mathbb{E}_{(s_t, g_t, a_t) \in \mathcal{D}_{\mathrm{BC}}} \log \pi_B(a_t | s_t, g_t; \theta) \qquad \text{Naive BC loss}$$

$$\mathcal{L}_{\mathrm{abc}} = -\mathbb{E}_{(s_t, g_t, a_t) \in \mathcal{D}_{\mathrm{BC}}} \left[ \mathrm{clip}\left( \frac{\pi_B(a_t | s_t, g_t; \theta)}{\pi_B(a_t | s_t, g_t; \theta_{\mathrm{old}})}, 1 - \epsilon, 1 + \epsilon \right) \right]$$

# Block Environment



- [1, 2] blocks
- State policy
  - Current position & rotation of blocks
  - Target position & rotation of blocks

# Evaluation: Skills to learn in the block environment

**Pushing**

**Flipping**

**Picking up**

**Stacking**



Transparent blocks mark the goal state.

# Generalize to unseen goals without manual curricula

- PPO (Schulman et al., 2017) baseline without curriculum fails to learn



| Push | Flip | Pick-and-place | Stack |

Success rate (%) vs Training steps (x100)

● No curriculum

# Generalize to unseen goals without manual curricula

- PPO (Schulman et al., 2017) baseline completely fails to learn
- Domain knowledge-based manual curriculum is insufficient

goal distance ratio
goal rotation weight
probability of pick-and-place
probability of stacking



Push      Flip      Pick-and-place      Stack

Success rate (%)

1 block
2 blocks
1 block
2 blocks

Training steps (x100)

● No curriculum     ● Curriculum

# Generalize to unseen goals without manual curricula

- PPO (Schulman et al., 2017) baseline completely fails to learn
- Domain knowledge-based manual curriculum is insufficient
- Asymmetric self-play zero-shot generalizes to all tasks

# Discovery of Novel Goals / Solutions



Novel Goals



Novel Solutions

# Discovery of Novel Goals / Solutions



Novel Goals

Novel Solutions

# Ablation: ABC is critical

- ⬤ ABC — Full setup with ABC
- 🟢 No ABC — No behavioral cloning loss in Bob's training

# Ablation: ABC is critical



No ABC



With ABC

# Ablation: ABC is critical

● ABC — Full setup with ABC

● No demonstration filter — Include all trajectories from Alice no matter Bob fails on this goal or not.



**Push**      **Flip**      **Pick-and-place**      **Stack**

2 blocks

1 block

Success rate (%)

Training steps (x100)

# Ablation: ABC is critical

● ABC — Full setup with ABC

● No BC loss clipping — No PPO-style loss clipping in ABC loss

# ShapeNet Environment



- [1, 10] objects from ShapeNet (Chang et al., 2015)
- Observation space: State + Vision
  - Current position & rotation of blocks
  - Target position & rotation of blocks
  - Images from front and wrist cameras
  - Target front-camera image

# Holdout tasks with unseen objects



| Table setting | Mini chess | Rainbow | Ball-capture | Tangram |

Initial State

Goal State

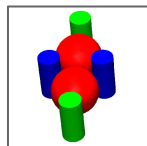# Zero-shot Generalization

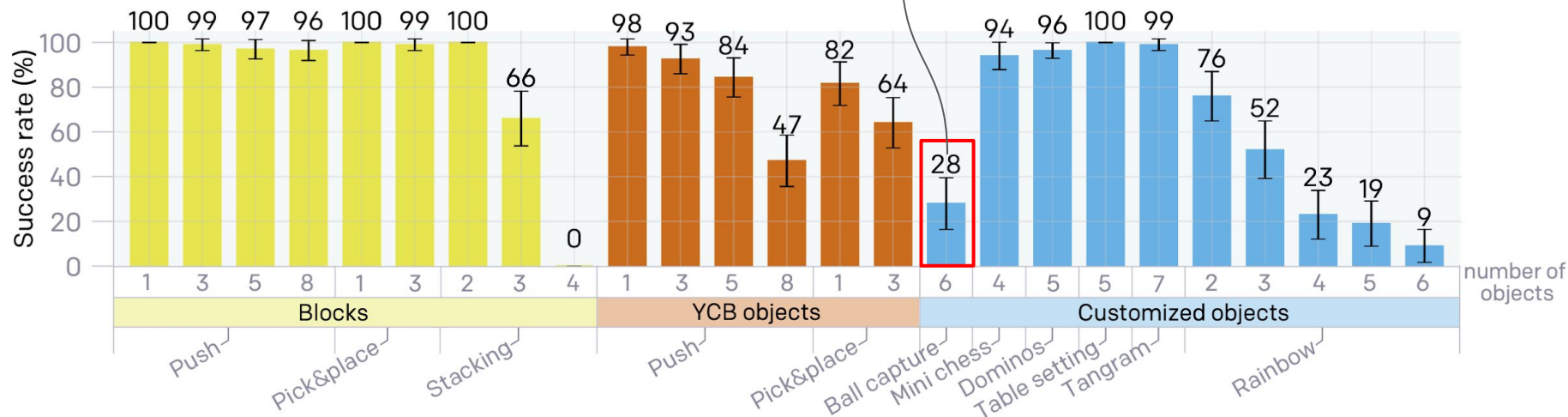Check out more videos at

https://robotics-self-play.github.io/

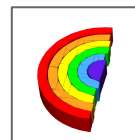# Zero-shot generalization

# Zero-shot generalization



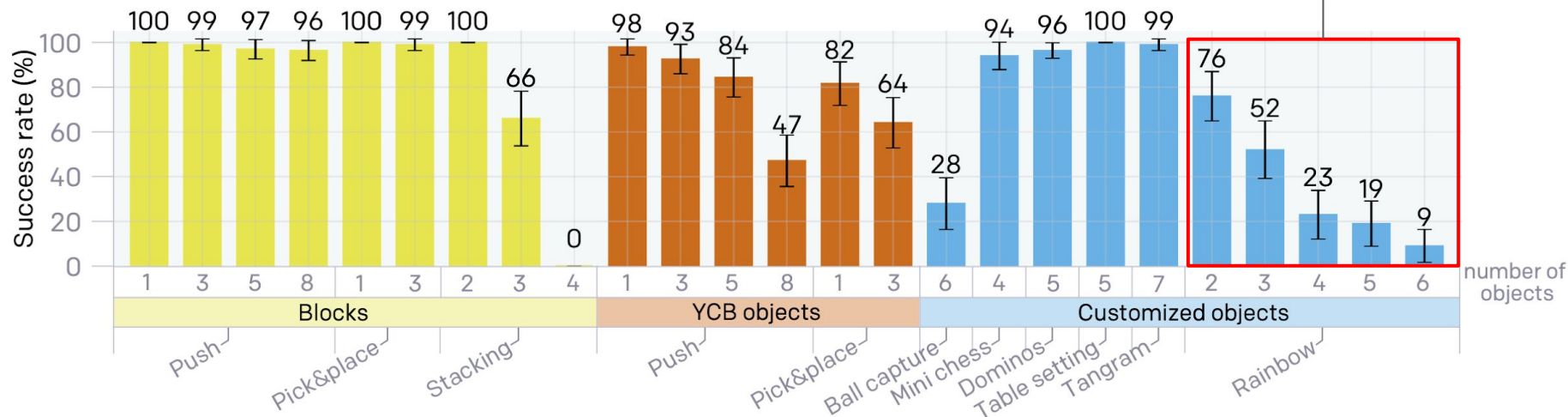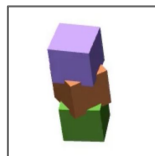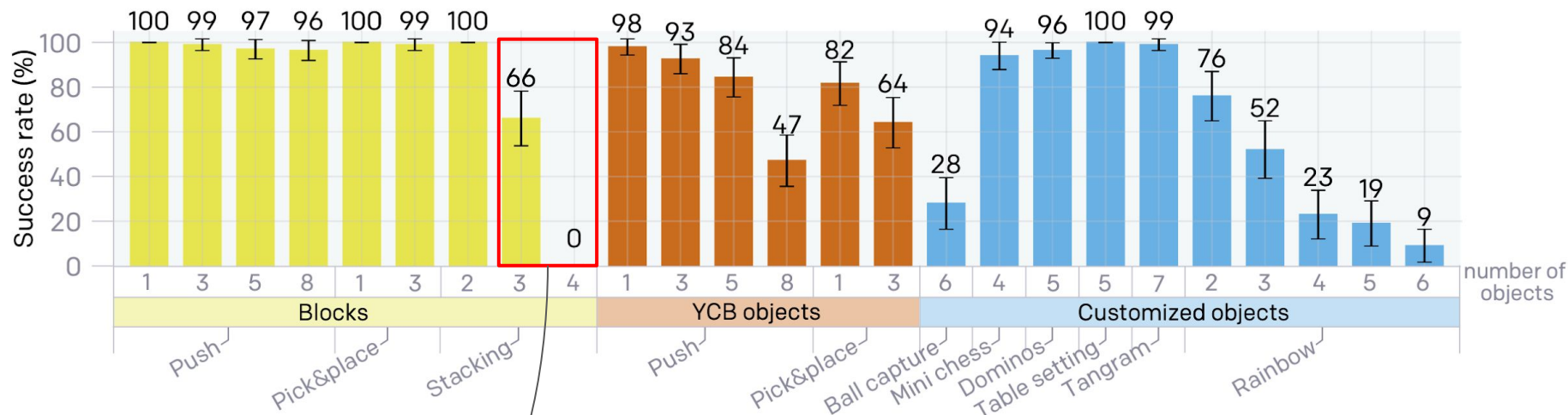Delicate handling of rolling objects and lifting skills

# Zero-shot generalization



Understand concave shapes

# Zero-shot generalization



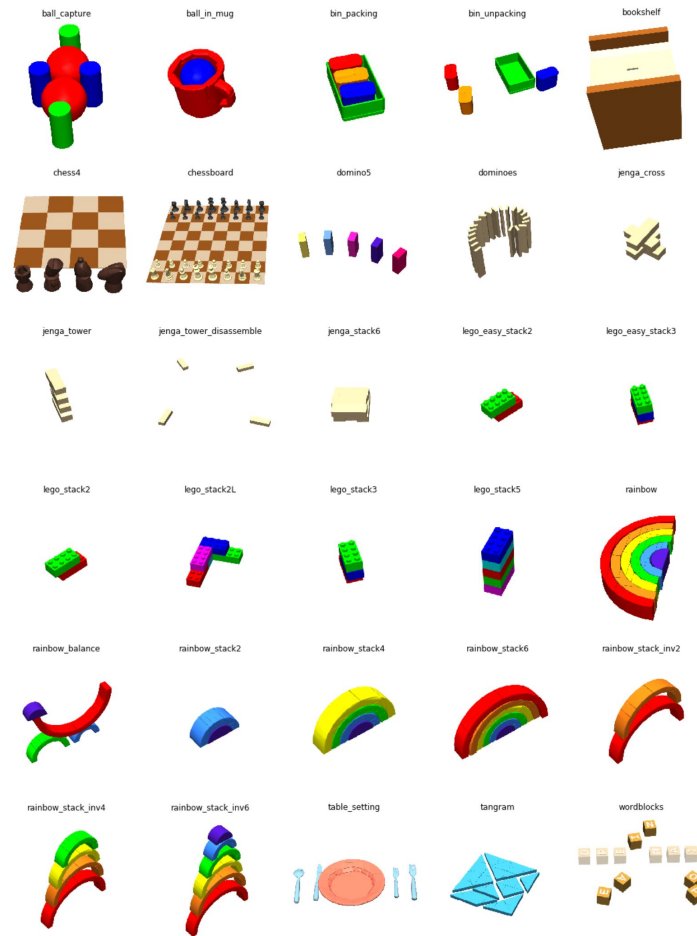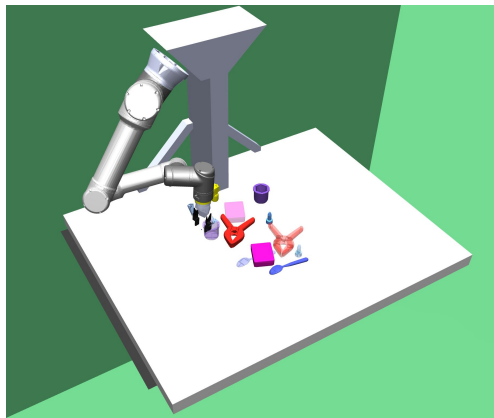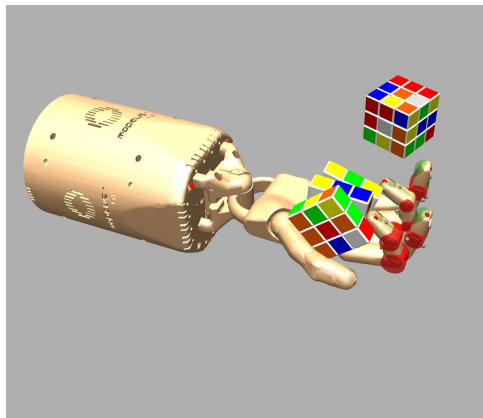Understand what's the correct order of placing objects.
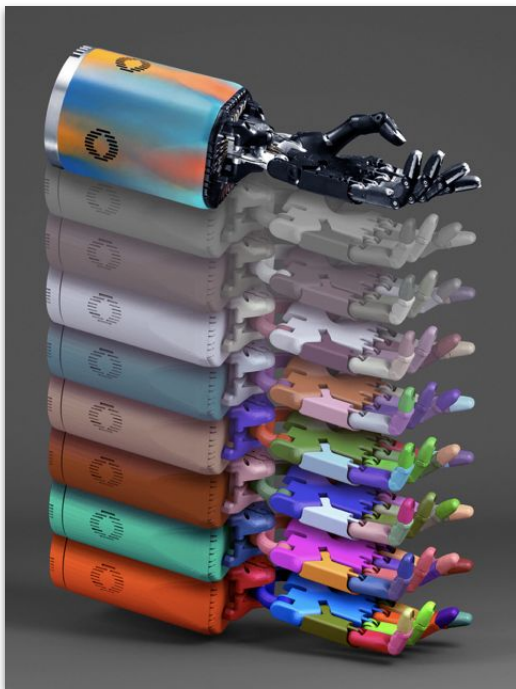
# Conclusion

Asymmetric self-play can:

1. Train a policy that can **zero-shot generalize** to many unseen robotic manipulation tasks.

2. **Alleviate** the importance of **manual curriculum**.

3. Alice Behavior Cloning (ABC) is crucial.

# Announce: `robogym`

https://github.com/openai/robogym

A simulation framework that uses OpenAI gym and MuJoCo simulator, including two environments: (1) in-hand manipulation with Rubik's cube; (2) table-top rearrange with one robot arm + gripper..
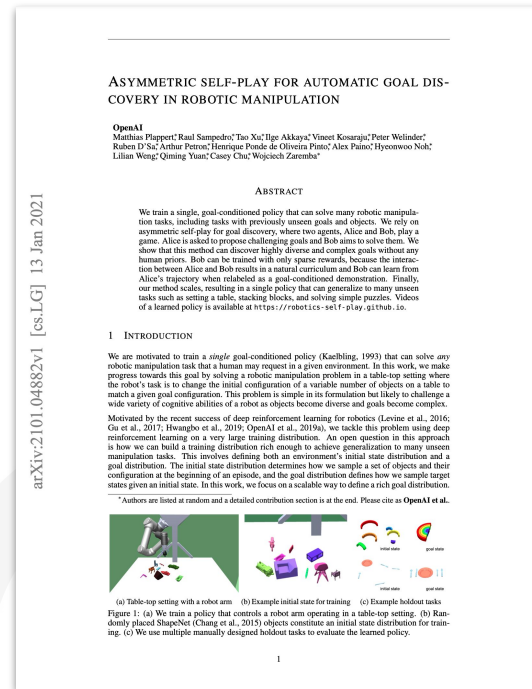
openai.com/blog/learning-dexterity

openai.com/blog/solving-rubiks-cube

arxiv.org/abs/2101.04882

# Thank you!

@lilianweng

🍑 lilianweng.github.io/lil-log