



# **Enhancing Language-Conditioned Robotic Manipulation through Prompt Tuning in a Missing Color Environment**

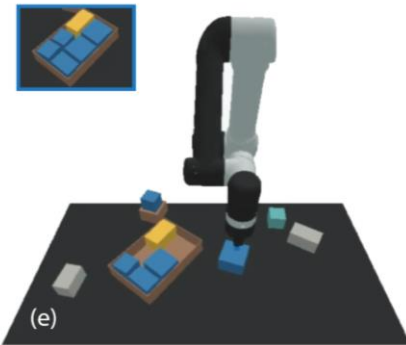
Baoming Zhang, Ouyang Xu

Presented on 12/04/2024

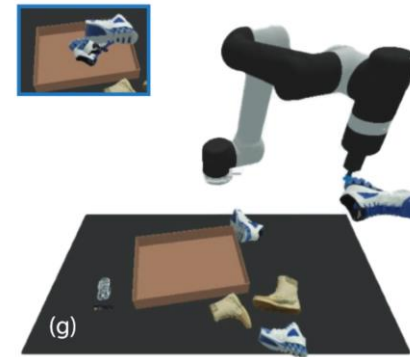
# Motivation

- Leverage large visual-language model, such as CLIP[1], on robotic grasping
- Robustness on text prompt corruptions: prompt engineering

## Language-conditioned tasks:





Pack all the yellow and blue blocks in the brown box



Pack all the blue and black sneaker objects in the brown box

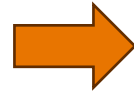
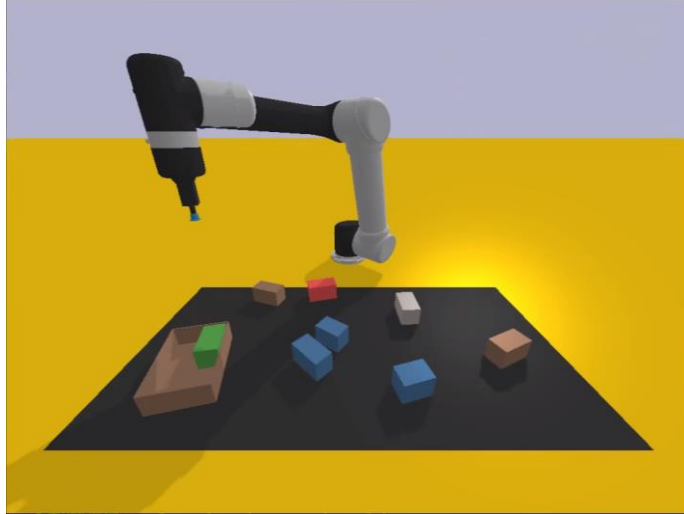
## Examples from CLIP:

Caltech101	Prompt	Accuracy
	a [CLASS].	82.68
	a photo of [CLASS].	80.81
	a photo of a [CLASS].	86.29

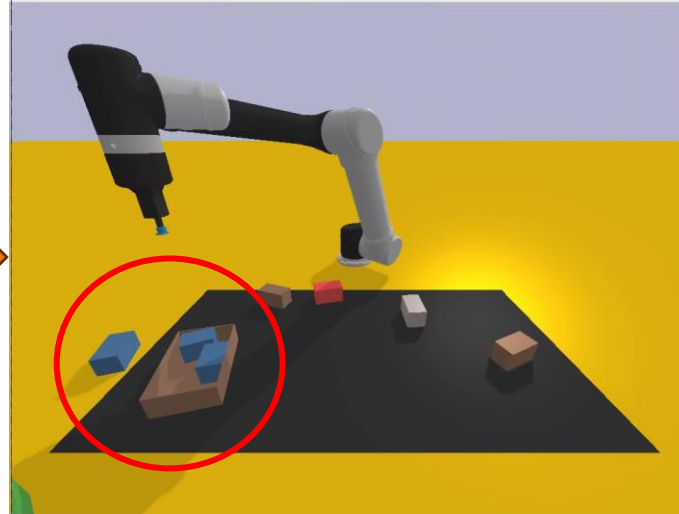
Flowers102	Prompt	Accuracy
	a photo of a [CLASS].	60.86
	a flower photo of a [CLASS].	65.81
	a photo of a [CLASS], a type of flower.	66.14

# Motivation

packing box pairs:  
all the blue and green blocks



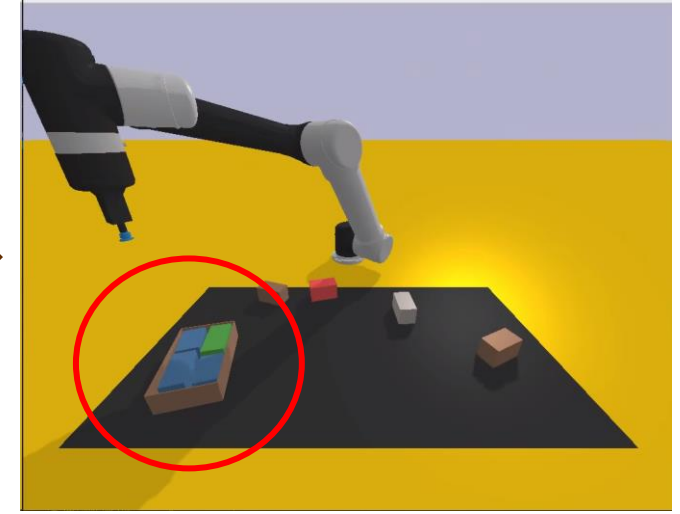
Template 1



Fail x



Template 3



Success ✓

Text Prompt	Task success scores (%)
1. 'pack all the [colors] blocks';	90.5
2. 'pack all the [colors] blocks into the box';	92.1
3. 'pack all the [colors] blocks into the brown box.'	97.1

# Problem Formulation (based on CLIPORT[1])

- **Objective:** Learn a goal-conditioned policy  $\pi$  that outputs actions  $\mathbf{a}_t$  based on inputs:

- $\gamma_t = (\mathbf{o}_t, \mathbf{l}_t)$ , where:

- \*  $\mathbf{o}_t$ : Visual observation.

- \*  $\mathbf{l}_t$ : English language instruction.

- **Policy Definition:**

$$\pi(\mathbf{o}_t, \mathbf{l}_t) \rightarrow \mathbf{a}_t = (\mathcal{T}_{\text{pick}}, \mathcal{T}_{\text{place}}) \in \mathcal{A}$$

- $\mathbf{a}_t$ : End-effector poses for:

- \*  $\mathcal{T}_{\text{pick}}$ : Picking.

- \*  $\mathcal{T}_{\text{place}}$ : Placing.

- **Task Focus:** Tabletop tasks where:

- $\mathcal{T}_{\text{pick}}, \mathcal{T}_{\text{place}} \in \mathbf{SE}(2)$ .

- **Visual Observations:**

- Top-down orthographic RGB-D reconstructions.

- Each pixel corresponds to a point in 3D space.

# Problem Formulation

- **Language Instructions:**

- Single goal descriptions:

- \* Example: *“Pack all the blue and yellow boxes in the brown box.”*

- **Dataset  $\mathcal{D}$ :**

- $n$  expert demonstrations:

$$\mathcal{D} = \{\zeta_1, \zeta_2, \dots, \zeta_n\}.$$

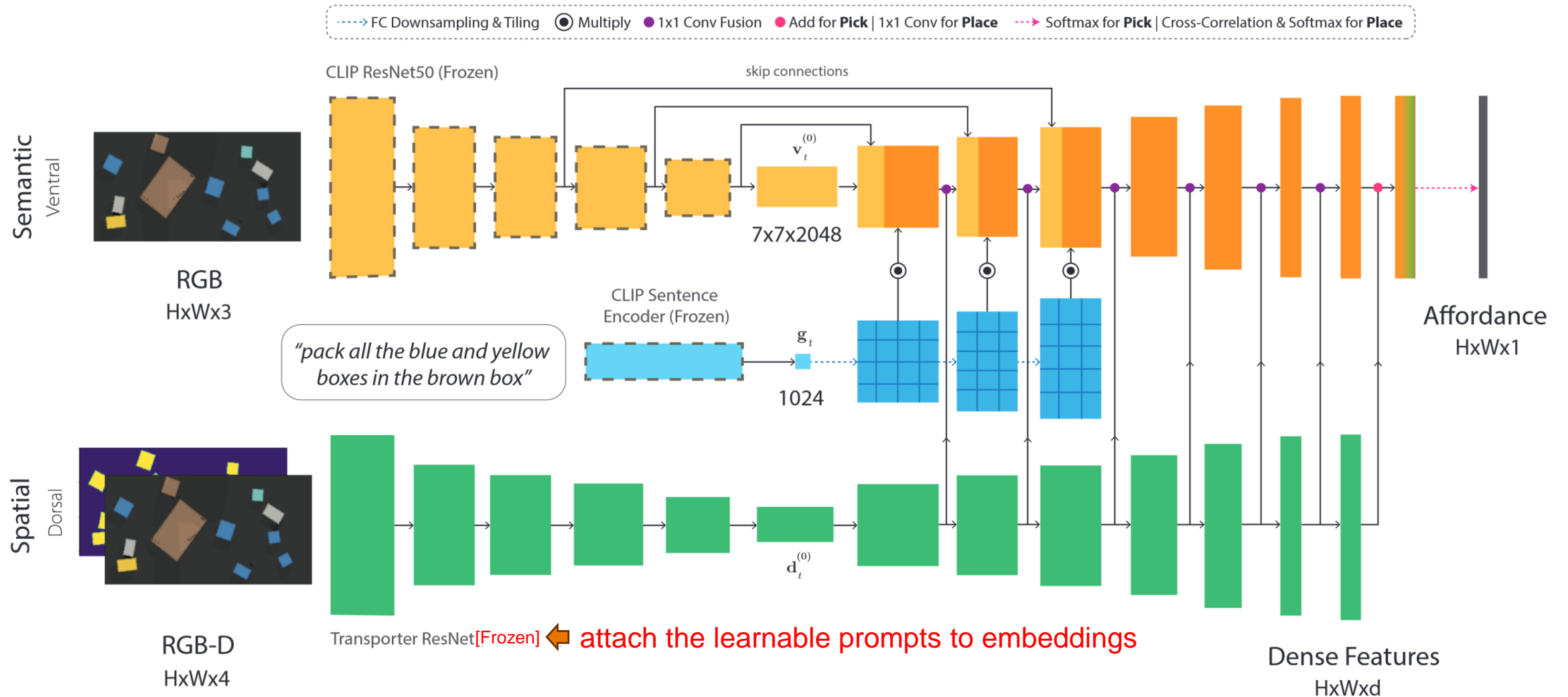
- Each demonstration  $\zeta_i$  contains input-action pairs:

$$\zeta_i = \{(\mathbf{o}_1, \mathbf{l}_1, \mathbf{a}_1), (\mathbf{o}_2, \mathbf{l}_2, \mathbf{a}_2), \dots\}.$$

- Actions  $\mathbf{a}_t = (\mathcal{T}_{\text{pick}}, \mathcal{T}_{\text{place}})$ : Expert pick-and-place coordinates.

- **Supervision:** Expert demonstrations are used to train the policy  $\pi$ .

# Foundational work



An overview of the semantic and spatial streams of foundation work CLIPORT[1].  
What we did for the missing color environment is in red.

# Transporter for Pick-and-Place

- **Overview:** Policy  $\pi$  is trained using Transporter [2] for spatial manipulation.
  - Two stages:
    1. Attend to a local region to determine the pick location.
    2. Compute the placement location using cross-correlation of deep visual features.
- **Policy Components:**
  - Two action-value modules (Q-functions):
    - \*  $Q_{\text{pick}}$ : Identifies the pick location.
    - \*  $Q_{\text{place}}$ : Determines the placement location conditioned on the pick action.
- **Place Module:**
  - Query FCN  $\Phi_{\text{query}}$  processes:
    - \*  $\gamma_t[\mathcal{T}_{\text{pick}}]$ :  $c \times c$  crop around  $\mathcal{T}_{\text{pick}}$ .
    - \*  $I_t$ : Language instruction.
  - Key FCN  $\Phi_{\text{key}}$  processes full input  $\gamma_t$ .
  - Placement action-values  $Q_{\text{place}}$ :

$$Q_{\text{place}}(\Delta\tau | \gamma_t, \mathcal{T}_{\text{pick}}) = (\Phi_{\text{query}}(\gamma_t[\mathcal{T}_{\text{pick}}]) * \Phi_{\text{key}}(\gamma_t))[\Delta\tau]$$

# Prompt Tuning

$$Q_{\text{place}}(\Delta\tau|\gamma_t, \mathcal{T}_{\text{pick}}) = (\Phi_{\text{query}}(\gamma_t[\mathcal{T}_{\text{pick}}]) * \Phi_{\text{key}}(\gamma_t))[\Delta\tau]$$

The prompts are appended as additional dimensions to the query and key embeddings

## Mathematical Formulation:

$$\Phi'_{\text{query}} = \text{concat}(\Phi_{\text{query}}(\gamma_t[\mathcal{T}_{\text{pick}}]), P_{\text{query}}, \text{dim} = -1)$$

$$\Phi'_{\text{key}} = \text{concat}(\Phi_{\text{key}}(\gamma_t), P_{\text{key}}, \text{dim} = -1)$$

## Dimensionality Impact:

- $P_{\text{query}} \in \mathbb{R}^{c \times c \times d_p}$ :  $d_p$  represents the prompt-specific channels.
- $P_{\text{key}} \in \mathbb{R}^{H \times W \times d_p}$ : Matches the spatial dimensions of the key embedding.

The new dimensions become:

$$\Phi'_{\text{query}} \in \mathbb{R}^{c \times c \times (d+d_p)}$$

$$\Phi'_{\text{key}} \in \mathbb{R}^{H \times W \times (d+d_p)}$$

Updated version:  $Q'_{\text{place}}(\Delta\tau|\gamma_t, \mathcal{T}_{\text{pick}}) = (\Phi'_{\text{query}}(\gamma_t[\mathcal{T}_{\text{pick}}]) * \Phi'_{\text{key}}(\gamma_t))[\Delta\tau]$



# Task Details

Task	multi-step sequencing	unseen colors	unseen objects	language instruction
1. packing-seen-google-objects-seq <sup>§</sup>	✓	✗	✗	step
2. packing-unseen-google-objects-seq <sup>§</sup>	✓	✓	✓	step
3. packing-seen-google-objects-group <sup>*§</sup>	✗	✗	✗	goal
4. packing-unseen-google-objects-group <sup>*§</sup>	✗	✓	✓	goal

<sup>§</sup>tasks that are commonly found in industry.

<sup>\*</sup>tasks that have more than one correct sequence of actions.

- **Selected Tasks:**

- 4 out of 10 language-conditioned tasks from the Ravens benchmark set [2].
- All tasks involve:
  - \* Precise placing.
  - \* Multimodal placing.

- **Language Templates for Training:**

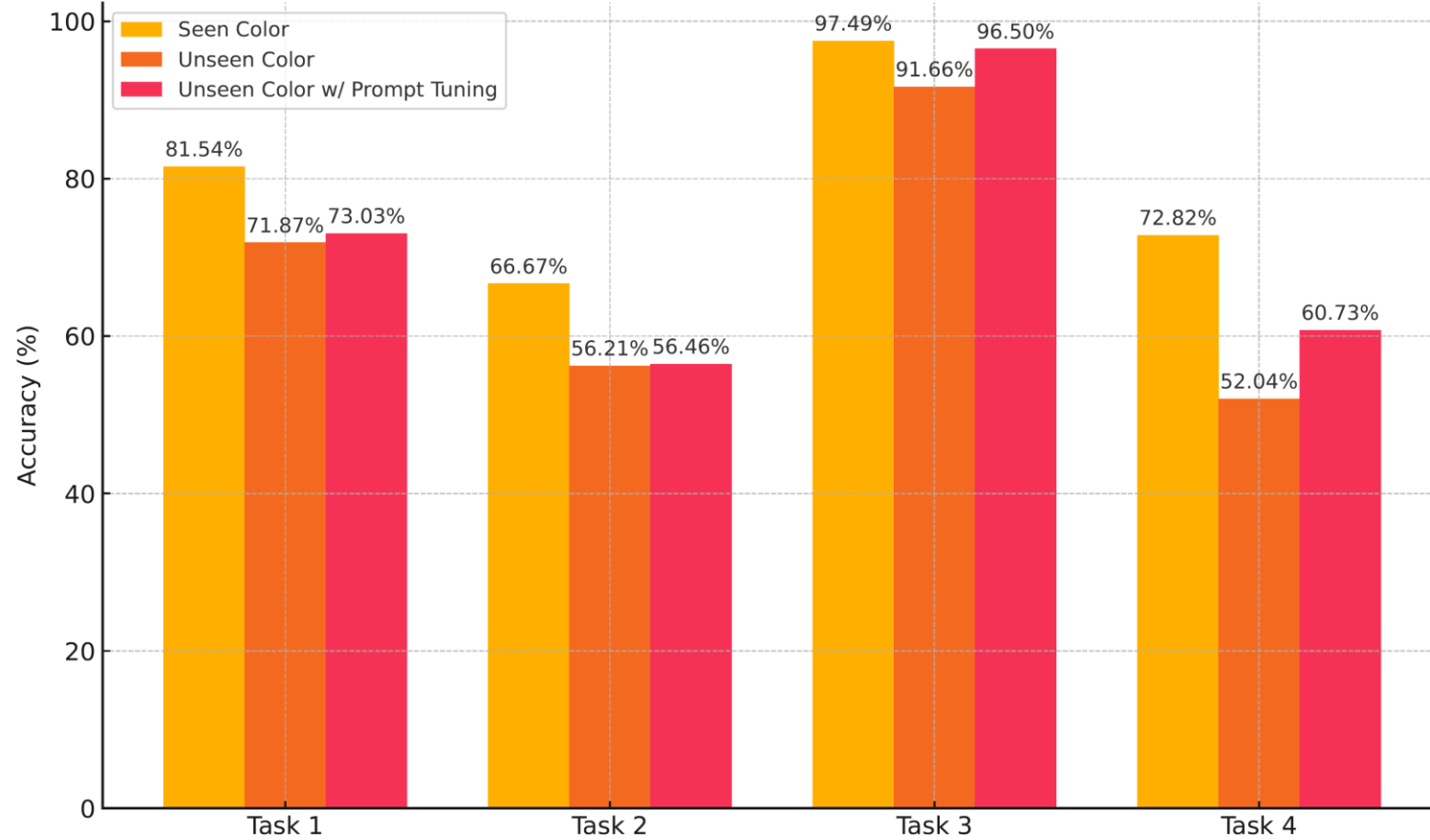
- Language instructions are distributed evenly across three templates:
  - \* Template 1: 'pack all the [colors] blocks'.
  - \* Template 2: 'pack all the [colors] blocks into the box'.
  - \* Template 3: 'pack all the [colors] blocks into the brown box'.

Other Training details:

Simulation environments (Ravens with PyBullet)  
The foundation model is frozen, only the prompts are trained

# Results

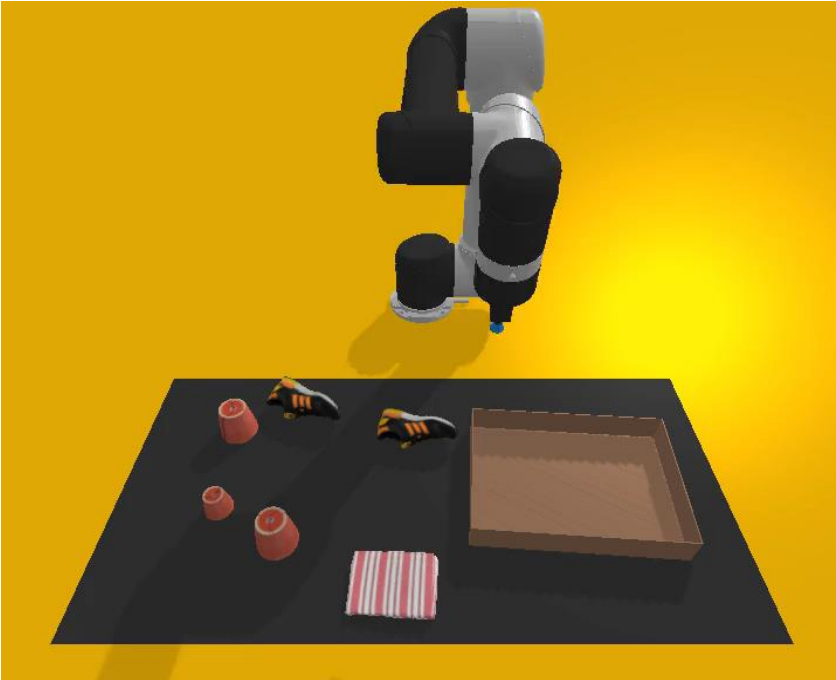
Task	Seen Color	Unseen Color	Unseen Color w/ Prompt Tuning
<b>Task1</b>	81.54%	71.87%	<b>73.03%</b>
<b>Task2</b>	66.67%	56.21%	<b>56.46%</b>
<b>Task3</b>	97.49%	91.66%	<b>96.50%</b>
<b>Task4</b>	72.82%	52.04%	<b>60.73%</b>



# Demo

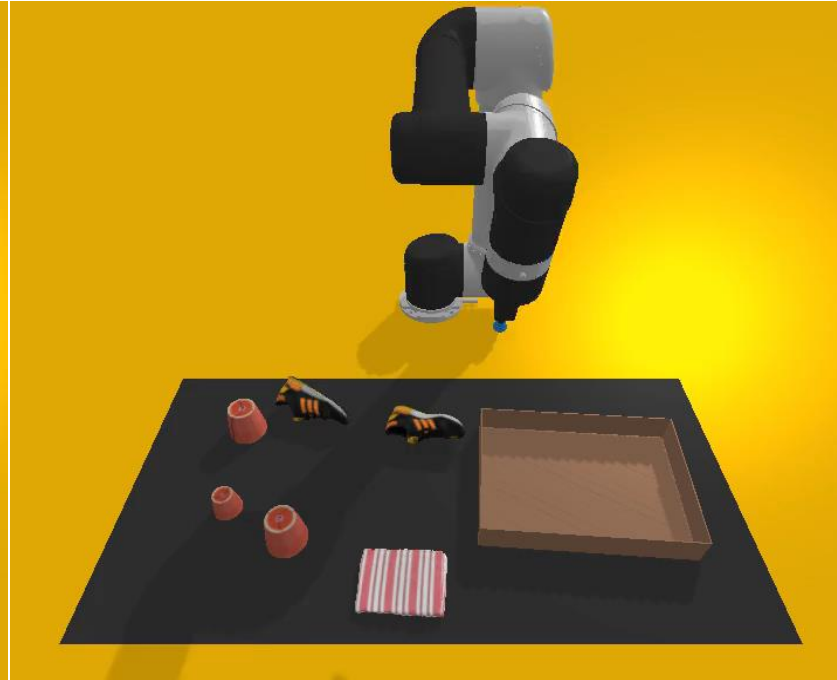
Task: Pack all the 'objects' in the brown box

Pack all the red cups  
in the brown box



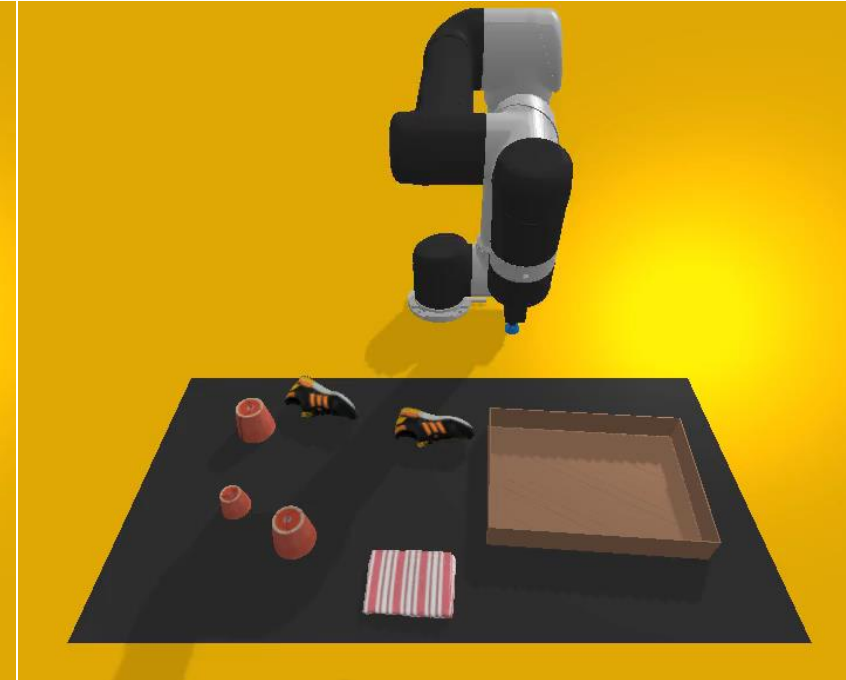
**Seen Color**

Pack all the red cups in the box



**Unseen Color**

Pack all the red cups in the box  
w/ Prompt Tuning



**Unseen Color w/ Prompt Tuning**



Thanks Everyone!