# Do As I Can, Not As I Say: Grounding Language in Robotic Affordance

He Zhi, CSE of SYSU

Email: hezh58@mail2.sysu.edu.cn

LLM has no experience, but has knowledge

Robot has no knowledge, but has experience

LLM **Say** knowledge

$+$

Robot **Can** do experience

$=$

# SayCan

SayCan is to ground large language models through value functions——
affordance functions that capture the log likelihood that a particular skill
will be able to succeed in the current state.

# Methodology

- i： Instruction

  How would you put an apple on the table?

- $\pi$： action in Robot action space

  $l_\pi$ ：action natural language description

  Find an apple

  Find a coke

  Pick up the apple

  Place the apple

  ......

- s： environment state



Gripper Height
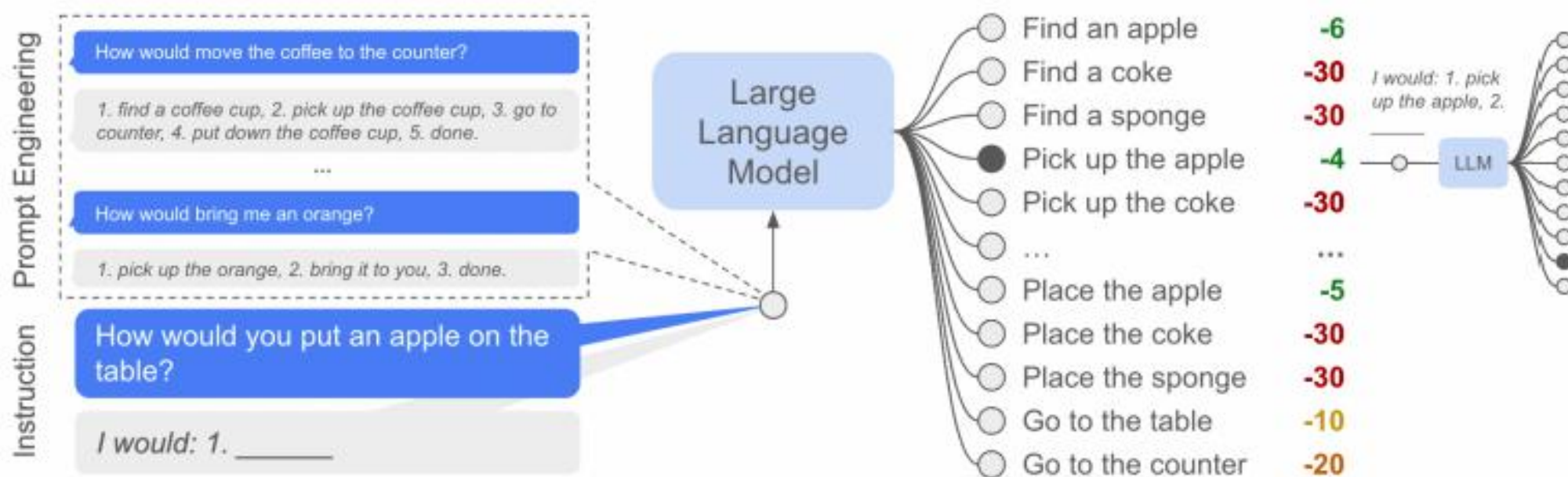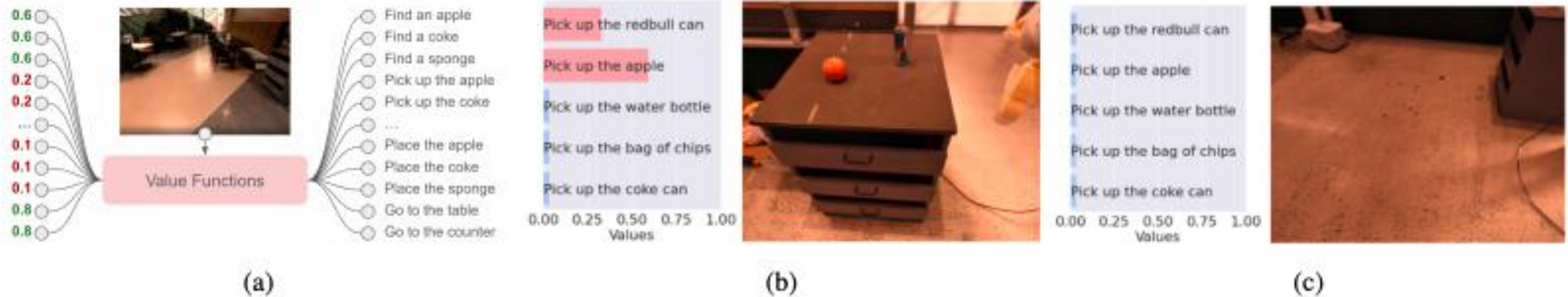
Rotation to Go

Closure to Go

......

+

# Methodology

LLM provides $p(l_\pi|i)$

# Methodology

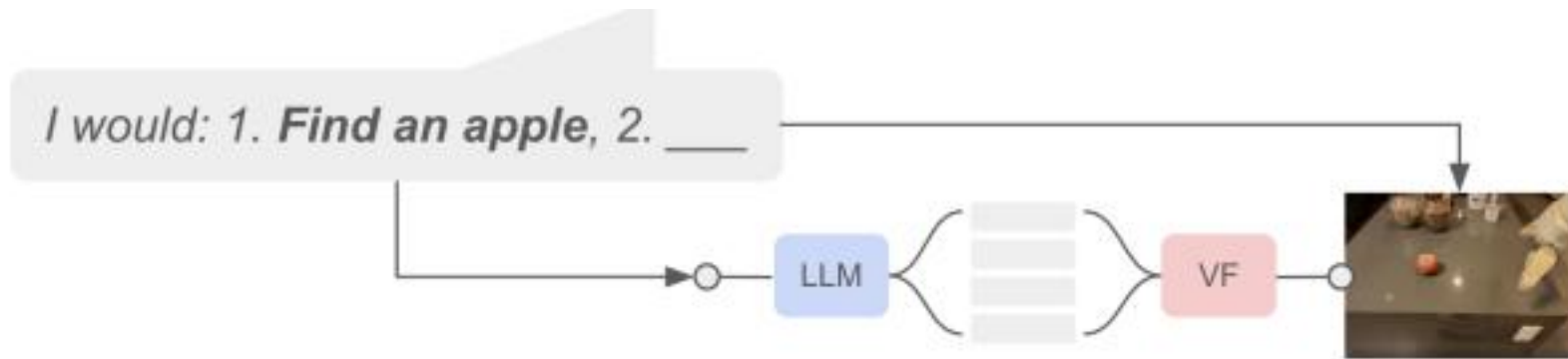Agent (Robot, etc) provides $p(c_\pi | l_\pi, s)$



(a)                                      (b)                                      (c)

# Methodology

$$p(c_\pi | i, s, l_\pi) \propto p(l_\pi | i) p(c_\pi | l_\pi, s)$$

$$\pi = \arg\max_{\pi \in \prod} p(l_\pi | i) p(c_\pi | l_\pi, s)$$



I would: 1. **Find an apple**, 2. ____

# Methodology

---

**Algorithm 1** SayCan

---

**Given:** A high level instruction $i$, state $s_0$, and a set of skills $\Pi$ and their language descriptions $\ell_\Pi$

1: $n = 0, \pi = \emptyset$
2: **while** $\ell_{\pi_{n-1}} \neq$ "done" **do**
3: $\quad \mathcal{C} = \emptyset$
4: $\quad$ **for** $\pi \in \Pi$ and $\ell_\pi \in \ell_\Pi$ **do**
5: $\quad\quad p_\pi^{\text{LLM}} = p(\ell_\pi | i, \ell_{\pi_{n-1}}, ..., \ell_{\pi_0})$ $\qquad\qquad\qquad\qquad\qquad$ ▷ Evaluate scoring of LLM
6: $\quad\quad p_\pi^{\text{affordance}} = p(c_\pi | s_n, \ell_\pi)$ $\qquad\qquad\qquad\qquad\qquad$ ▷ Evaluate affordance function
7: $\quad\quad p_\pi^{\text{combined}} = p_\pi^{\text{affordance}} p_\pi^{\text{LLM}}$
8: $\quad\quad \mathcal{C} = \mathcal{C} \cup p_\pi^{\text{combined}}$
9: $\quad$ **end for**
10: $\quad \pi_n = \arg\max_{\pi \in \Pi} \mathcal{C}$
11: $\quad$ Execute $\pi_n(s_n)$ in the environment, updating state $s_{n+1}$
12: $\quad n = n + 1$
13: **end while**

---

# Model architecture——Agent



Figure 9: Network architecture in RL policy

《Mt-opt: Continuous multi-task robotic reinforcement learning at scale》

# Model architecture——Agent



Figure 10: Network architecture in BC policy

《Bc-z:Zero-shot task generalization with robotic imitation learning》

# Model architecture——Agent

expert DM data: $\{T_1, T_2, T_3, \ldots T_n\}$

$T_i = \langle s_1^i, a_1^i, s_2^i, a_2^i, \ldots s_T^i, a_T^i \rangle$

$\Downarrow$

$D = \{(s_1^i, a_1^i), (s_2^i, a_2^i), \ldots \}$

$\Downarrow$

classifier



Figure 10: Network architecture in BC policy

《Bc-z：Zero-shot task generalization with robotic imitation learning》

# Model architecture——LLM

- PaLM 《PaLM: Scaling Language Modeling with Pathways》

- FLAN 《FINETUNED LANGUAGE MODELS ARE ZERO-SHOT LEARNERS》

# Experiment／Result



(a) "I just worked out, can you bring me a drink and a snack to recover?"



Task: **move all the blocks into their matching colored bowls**.
Step 1. pick up the blue block and place it in the blue bowl
Step 2. pick up the green block and place it in the green bowl
Step 3. pick up the yellow block and place it in the yellow bowl

ViLD-based Affordances     Initial State     Final State

# Experiment/Result

| Instruction Family | Num | Explanation | Example Instruction |
|---|---|---|---|
| NL Single Primitive | 15 | NL queries for a single primitive | Let go of the coke can |
| NL Nouns | 15 | NL queries focused on abstract nouns | Bring me a fruit |
| NL Verbs | 15 | NL queries focused on abstract verbs | Restock the rice chips on the far counter |
| Structured Language | 15 | Structured language queries, mirror NL Verbs | Move the rice chips to the far counter. |
| Embodiment | 11 | Queries to test SayCan's understanding of the current state of the environment and robot | Put the coke on the counter. (starting from different completion stages) |
| Crowd-Sourced | 15 | Queries in unstructured formats | My favorite drink is redbull, bring one |
| Long-Horizon | 15 | Long-horizon queries that require many steps of reasoning | I spilled my coke on the table, throw it away and bring me something to clean |

| Instruction |
|---|
| How would you bring me lime drink |
| How would you bring me something to clean the kitchen with |
| How would you bring me something to eat |
| How would you put the grapefruit drink on the close counter |
| How would you move the sugary drink to the far counter |
| How would you move something with caffine from the table to the close counter |
| How would you bring me an energy bar |
| How would you bring me something to quench my thirst |
| How would you bring me a fruit |
| How would you bring me a fruit from the close counter |
| How would you bring me something that is not a fruit from the close counter |
| How would you bring me a soda from the table |
| How would you bring me a soda |
| How would you bring me a bag of chips from close counter |
| How would you bring me a snack |

(c) NL Nouns

| Instruction |
|---|
| I opened a pepsi earlier. How would you bring me an open can? |
| I spilled my coke, can you bring me a replacement? |
| I spilled my coke, can you bring me something to clean it up? |
| I accidentally dropped that jalapeno chip bag after eating it. Would you mind throwing it away? |
| I like fruits, can you bring me something I'd like? |
| There is a close counter, far counter, and table. How would you visit all the locations? |
| There is a close counter, trash can, and table. How would you visit all the locations? |
| Redbull is my favorite drink, can I have one please? |
| Would you bring me a coke can? |
| Please, move the pepsi to the close counter |
| Please, move the ppsi(intentional typo) to the close cuonter |
| Can you move the coke can to the far counter? |
| Can you move coke can to far counter? |
| Would you throw away the bag of chips for me? |
| Would you throw away the bag of chpis(intentional typo) for me? |

(f) Crowd-Sourced

# Experiment/Result

- plan success rate

- execution success rate

Handwritten annotations: *no value function*, *→ Generative*

| | | Mock Kitchen | | Kitchen | | No Affordance | | No LLM | |
|---|---|---|---|---|---|---|---|---|---|
| | | PaLM-SayCan | PaLM-SayCan | PaLM-SayCan | PaLM-SayCan | No VF | Gen. | BC NL | BC USE |
| **Family** | **Num** | Plan | Execute | Plan | Execute | Plan | Plan | Execute | Execute |
| NL Single | 15 | 100% | 100% | 93% | 87% | 73% | 87% | 0% | 60% |
| NL Nouns | 15 | 67% | 47% | 60% | 40% | 53% | 53% | 0% | 0% |
| NL Verbs | 15 | 100% | 93% | 93% | 73% | 87% | 93% | 0% | 0% |
| Structured | 15 | 93% | 87% | 93% | 47% | 93% | 100% | 0% | 0% |
| Embodiment | 11 | 64% | 55% | 64% | 55% | 18% | 36% | 0% | 0% |
| Crowd Sourced | 15 | 87% | 87% | 73% | 60% | 67% | 80% | 0% | 0% |
| Long-Horizon | 15 | 73% | 47% | 73% | 47% | 67% | 60% | 0% | 0% |
| Total | 101 | 84% | 74% | 81% | 60% | 67% | 74% | 0% | 9% |

Table 2: Success rates of integration by family. PaLM-SayCan achieves a planning success rate of 84% and

# What to do next?

Sparse reward!

No middle-step reward!

# What to do next?

Sparse reward!

No middle-step reward!

human: How would you put an apple on the table?
robot: ① Go to the counter → Find the apple → Pick up the apple → Go back to the table
       ↳ Place the apple ⟹ reward : +1

② Go to the counter → Find the apple → Pick up the apple → Pick up the apple
   Pick up the apple → Pick up the apple → …… → Pick up the Coak ⟹ reward : 0

# Code/Dataset

- https://github.com/google-research/google-research/tree/master/saycan

- https://github.com/say-can/say-can.github.io/tree/main/data