# A Survey of Robotic Language Grounding:
## Tradeoffs between Symbols and Embeddings

Vanya Cohen, **Jason Xinyu Liu**, Raymond Mooney, Stefanie Tellex, David Watkins

IJCAI 2024 Survey Track

August 8th, 2024

# Robotic Language Grounding

Connect linguistic elements in language to the robot's perception of and actions in the physical world.

Liu et al. IJCAI 2024 Survey Track. A Survey of Robotic Language Grounding: Tradeoffs between Symbols and Embeddings

1

# Robotic Language Grounding

Connect linguistic elements in language to the robot's perception of and actions in the physical world.

1. What grounding representation to use?

# Robotic Language Grounding

Connect linguistic elements in language to the robot's perception of and actions in the physical world.

1. What grounding representation to use?

2. How to ground natural language to the grounding representation of choice?

# Robotic Language Grounding

**Input Command**

**Target Representation**

# Grounding Language to Symbols



**Input Command**

"Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot."

**Target Representation**

First Order Logic | Linear Temporal Logic

in(the(
λx.agent(x)),
the(λy.green(y) ∧
room(y))) | ♦(walmart) ∧
■¬(parking_lot)

**Discrete**
**More Structure**
**More Bias**

LEFT
Hsu et al., 2023

## Symbols

- Discrete

- More Structure; More bias

- Unambiguous

- Verifiable

- Interpretable

# Grounding Language to Embeddings



**Symbols**

- Discrete

- More Structure; More bias

- Unambiguous

- Verifiable

- Interpretable

**High-dimensional Embeddings**

- Continuous

- Less structure; More variance

- Adaptive

# Grounding Language to Logic

# Grounding Language to Logic



Pros

- Unambiguous semantics

- Verifiable

- Interpretable

- Reduce search space

# Grounding Language to Logic



**Pros**

- Unambiguous semantics
- Verifiable
- Interpretable
- Reduce search space

**Cons**

- Require manually defined structures
- Difficult to represent low-level control

# Grounding Language to Logic: Lang2LTL



Lang2LTL

- Natural language navigation command

- Modular system produces a grounded linear temporal logic (LTL) formula

- Given MDP definition

- Planner outputs a trajectory

# Grounding Language to Logic



**More Papers**

- Lang2LTL-2: Grounding Spatiotemporal Navigation Commands Using Large Language and Vision-Language Models [Liu et al. 2024]

- AutoTAMP: Autoregressive Task and Motion Planning with LLMs as Translators and Checkers [Chen et al. 2024]

- NL2TL: Transforming Natural Languages to Temporal Logics using Large Language Models [Chen et al. 2023]

- NL2LTL: a Python Package for Converting Natural Language (NL) Instructions to Linear Temporal Logic (LTL) Formulas [Fuggitti and Chakraborti 2023]

# Grounding Language to PDDL

# Grounding Language to PDDL

| | | |
|---|---|---|
| **Input Command** | | |
| "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." |
| **Target Representation** | | |
| First Order Logic | Linear Temporal Logic | PDDL |
| in(the( λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) |

**Discrete**
**More Structure**
**More Bias**

| LEFT | Lang2LTL | LLM+P |
|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 |

## Pros

- Sound
- Complete
- (Often) Optimal

# Grounding Language to PDDL



## Pros

- Sound

- Complete

- (Often) Optimal

## Cons

- Require manually defined structures

# Grounding Language to PDDL: LLM+P



LLM+P

- Natural language description of a planning problem

- LLM translates it to PDDL problem

- Given a PDDL domain description, i.e., action preconditions and effects

- Symbolic planner solves PDDL

# Grounding Language to PDDL



| | | | |
|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." |
| **Target Representation** | First Order Logic<br><br>in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | Linear Temporal Logic<br><br>◆(walmart) ∧<br>■¬(parking_lot) | PDDL<br><br>Initial:<br>  (ontable laptop)<br>  (on novel laptop)<br>Goal:<br>  (ontable novel) |

**Discrete**
**More Structure**
**More Bias**

LEFT
Hsu et al., 2023

Lang2LTL
Liu et al., 2023b

LLM+P
Liu et al., 2023

## More Papers

- Translating Natural Language to Planning Goals with Large-Language Models [Xie et al. 2023]
- Structured, Flexible, and Robust: Benchmarking and Improving Large Language Models Towards More Human-like Behavior in Out-of-distribution Reasoning Tasks [Collins et al. 23]
- Leveraging Pre-trained Large Language Models to Construct and Utilize World Models for Model-based Task Planning [Guan et al. 2023]
- PlanBench: An Extensible Benchmark for Evaluating Large Language Models on Planning and Reasoning about Change [Valmeekam et al. 2023]
- On the Planning Abilities of Large Language Models : A Critical Investigation [Valmeekam et al. 2023]

# Grounding Language to Code



| | | | |
|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." |

Target Representation:

First Order Logic
```
in(the(
λx.agent(x)),
the(λy.green(y) ∧
room(y)))
```

Linear Temporal Logic
```
♦(walmart) ∧
■¬(parking_lot)
```

PDDL
```
Initial:
    (ontable laptop)
    (on novel
    laptop)
Goal:
    (ontable novel)
```

Code: Python, JS...
```
big_block =
max(blocks,
key=lambda x:
x.volume())
pick_up(big_block)
```

Discrete
More Structure
More Bias

| LEFT | Lang2LTL | LLM+P | Code as Policies |
|---|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 | Liang et al., 2022 |

# Grounding Language to Code



| Input Command | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." |
|---|---|---|---|---|
| Target Representation | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... |
| | in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) |

Discrete
More Structure
More Bias

| LEFT | Lang2LTL | LLM+P | Code as Policies |
|---|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 | Liang et al., 2022 |

## Pros

- Flexible

- High-level plan and low-level control

# Grounding Language to Code



| Input Command | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." |
|---|---|---|---|---|
| Target Representation | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... |
| | in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) |

**Discrete More Structure More Bias** →

| LEFT Hsu et al., 2023 | Lang2LTL Liu et al., 2023b | LLM+P Liu et al., 2023 | Code as Policies Liang et al., 2022 |

## Pros

• Flexible

• High-level plan and low-level control

## Cons

• Require predefined perception and control models in specific domains

# Grounding Language to Code: Code as Policies



Code as Policies

- Natural language command
- Given predefined perception and control models
- Code-writing LLM outputs executable code

# Grounding Language to Code



| | | | |
|---|---|---|---|
| **Input Command** | | | |
| "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." |
| **Target Representation** | | | |
| First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... |
| in(the( λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) |
| **LEFT** Hsu et al., 2023 | **Lang2LTL** Liu et al., 2023b | **LLM+P** Liu et al., 2023 | **Code as Policies** Liang et al., 2022 |

**Discrete**
**More Structure**
**More Bias**

## More Papers

- Embodied AI with Two Arms: Zero-shot Learning, Safety and Modularity [Varley et al. 2024]
- ProgPrompt: Generating Situated Robot Task Plans using Large Language Models [Singh et al. 2023]
- Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language [Zeng et al. 2023]
- ITP: Interactive Task Planning with Language Models [Li et al. 2023]
- Voyager: An Open-ended Embodied Agent with Large Language Models [Wang et al. 2023]

# Grounding Language to Predefined Skills



| | | | | |
|---|---|---|---|---|
| **Input Command** | | | | |
| "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" |
| **Target Representation** | | | | |
| First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills |
| in(the($\lambda$x.agent(x)), the($\lambda$y.green(y) $\wedge$ room(y))) | $\blacklozenge$(walmart) $\wedge$ $\blacksquare\neg$(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you |
| **LEFT** Hsu et al., 2023 | **Lang2LTL** Liu et al., 2023b | **LLM+P** Liu et al., 2023 | **Code as Policies** Liang et al., 2022 | **SayCan** Ahn et al. 2022 |

**Discrete**
**More Structure**
**More Bias**

# Grounding Language to Predefined Skills



| | | | | |
|---|---|---|---|---|
| **Input Command** | | | | |
| "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" |
| **Target Representation** | | | | |
| First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills |
| in(the( $\lambda$x.agent(x)), the($\lambda$y.green(y) $\wedge$ room(y))) | $\blacklozenge$(walmart) $\wedge$ $\blacksquare\neg$(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you |
| **LEFT** Hsu et al., 2023 | **Lang2LTL** Liu et al., 2023b | **LLM+P** Liu et al., 2023 | **Code as Policies** Liang et al., 2022 | **SayCan** Ahn et al. 2022 |

**Discrete**
**More Structure**
**More Bias**

Pros

- Adaptive

# Grounding Language to Predefined Skills



| | | | | |
|---|---|---|---|---|
| **Input Command** | | | | |
| "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" |
| **Target Representation** | | | | |
| First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills |
| in(the( λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you |
| **LEFT** Hsu et al., 2023 | **Lang2LTL** Liu et al., 2023b | **LLM+P** Liu et al., 2023 | **Code as Policies** Liang et al., 2022 | **SayCan** Ahn et al. 2022 |

**Discrete / More Structure / More Bias** →

## Pros

- Adaptive

## Cons

- Require predefined skills

- Possibly incorrect plans

# Grounding Language to Predefined Skills: SayCan



SayCan

- Natural language command

- LLM proposes candidate skills every step

- Pretrained value functions to rank available skills

- Language-conditioned policies execute the top skill

# Grounding Language to Predefined Skills



**More Papers**

- CAPE: Planning with Large Language Models via Corrective Re-prompting [Raman et al. 2024]

- Inner Monologue: Embodied Reasoning through Planning with Language Models [Huang et al. 2022]

- Language Models as Zero-shot Planners: Extracting Actionable Knowledge for Embodied Agent [Huang et al. 2022]

# Grounding Language to Subgoals



| | | | | | |
|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-a |

First Order Logic / Linear Temporal Logic / PDDL / Code: Python, JS... / Predefined Skills / Natural Language / Vision

**Target Representation:**

First Order Logic:
in(the(
λx.agent(x)),
the(λy.green(y) ∧
room(y)))

Linear Temporal Logic:
♦(walmart) ∧
■¬(parking_lot)

PDDL:
Initial:
(ontable laptop)
(on novel laptop)
Goal:
(ontable novel)

Code: Python, JS...:
big_block =
max(blocks,
key=lambda x:
x.volume())
pick_up(big_block)

Predefined Skills:
1. find a sponge
2. pick up the sponge
3. bring it to you

Natural Language:
Steps:
Step 1: chop tree block
Step 2: craft wooden stick...

Vision: Images

**Discrete**
**More Structure**
**More Bias**

| LEFT | Lang2LTL | LLM+P | Code as Policies | SayCan | VLP |
|---|---|---|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 | Liang et al., 2022 | Ahn et al. 2022 | Du et al., 2023 |

# Grounding Language to Subgoals



| | | | | | | |
|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-a |
| **Target Representation** | First Order Logic<br><br>in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | Linear Temporal Logic<br><br>♦(walmart) ∧<br>■¬(parking_lot) | PDDL<br><br>Initial:<br>(ontable laptop)<br>(on novel laptop)<br>Goal:<br>(ontable novel) | Code: Python, JS...<br><br>big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | Predefined Skills<br><br>1. find a sponge<br>2. pick up the sponge<br>3. bring it to you | Natural Language / Vision<br>Steps: / Images<br>Step 1: chop tree block<br>Step 2: craft wooden stick... |

**Discrete**
**More Structure**
**More Bias**

| LEFT | Lang2LTL | LLM+P | Code as Policies | SayCan | VLP |
|---|---|---|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 | Liang et al., 2022 | Ahn et al. 2022 | Du et al., 2023 |

## Pros

- Adaptive

# Grounding Language to Subgoals

| | | | | | |
|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-a |
| **Target Representation** | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills | Natural Language / Vision |
| | in(the( λx.agent(x)), the(λy.green(y) ∧ room(y))) | ◆(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you | Steps: Step 1: chop tree block Step 2: craft wooden stick... / Images |

**Discrete
More Structure
More Bias**

| LEFT | Lang2LTL | LLM+P | Code as Policies | SayCan | VLP |
|---|---|---|---|---|---|
| Hsu et al., 2023 | Liu et al., 2023b | Liu et al., 2023 | Liang et al., 2022 | Ahn et al. 2022 | Du et al., 2023 |

**Pros**

- Adaptive

**Cons**

- Require predefined skills
- Possibly incorrect plans

# Grounding Language to Subgoals: VLP



Video Language Planning (VLP)

- Tree search

- VLM proposes language subgoals

- Video model conditioned on text generates image subgoals

- Policy conditioned on image executes the plan

# Grounding Language to Subgoals



| | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills | Natural Language | Vision |
|---|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-a..." | |
| **Target Representation** | in(the( λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you | Steps: Step 1: chop tree block Step 2: craft wooden stick... | Images |

**Discrete / More Structure / More Bias** →

- LEFT — Hsu et al., 2023
- Lang2LTL — Liu et al., 2023b
- LLM+P — Liu et al., 2023
- Code as Policies — Liang et al., 2022
- SayCan — Ahn et al. 2022
- VLP — Du et al., 2023

## More Papers

- Zero-Shot Robotic Manipulation with Pretrained Image-Editing Diffusion Models [Black et al. 2023]

- UniSim: A Neural Closed-Loop Sensor Simulator [Yang et al. 2023]

- GAIA-1: A Generative World Model for Autonomous Driving [Hu et al. 2023]

# Grounding Language to Embeddings



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." | |
| **Target Representation** | First Order Logic<br>in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | Linear Temporal Logic<br>♦(walmart) ∧ ■¬(parking_lot) | PDDL<br>Initial:<br>(ontable laptop)<br>(on novel laptop)<br>Goal:<br>(ontable novel) | Code: Python, JS...<br>big_block = max(blocks, key=lambda x: x.volume())<br>pick_up(big_block) | Predefined Skills<br>1. find a sponge<br>2. pick up the sponge<br>3. bring it to you | Natural Language<br>Steps:<br>Step 1: chop tree block<br>Step 2: craft wooden stick... | Vision<br>Images | Vector<br>Model Activations |

**Discrete More Structure More Bias** ← → **Continuous Less Structure More Variance**

| LEFT<br>Hsu et al., 2023 | Lang2LTL<br>Liu et al., 2023b | LLM+P<br>Liu et al., 2023 | Code as Policies<br>Liang et al., 2022 | SayCan<br>Ahn et al. 2022 | VLP<br>Du et al., 2023 | VIMA<br>Jiang et al., 2023 |
|---|---|---|---|---|---|---|

# Grounding Language to Embeddings

| | Input Command | Target Representation | | |
|---|---|---|---|---|
| **First Order Logic** | "Go to the green room." | in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | | |
| **Linear Temporal Logic** | "Go to the store on Main Street, and always avoid the parking lot." | ♦(walmart) ∧ ■¬(parking_lot) | | |
| **PDDL** | "Move the novel onto the table." | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | | |
| **Code: Python, JS...** | "Pick up the largest block." | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | | |
| **Predefined Skills** | "I spilled my drink, can you help?" | 1. find a sponge 2. pick up the sponge 3. bring it to you | | |
| **Natural Language** | "Craft a diamond pick-axe." | Steps: Step 1: chop tree block Step 2: craft wooden stick... | | |
| **Vision** | | Images | | |
| **Vector** | | Model Activations | | |

Discrete More Structure More Bias → Continuous Less Structure More Variance

- **LEFT** Hsu et al., 2023
- **Lang2LTL** Liu et al., 2023b
- **LLM+P** Liu et al., 2023
- **Code as Policies** Liang et al., 2022
- **SayCan** Ahn et al. 2022
- **VLP** Du et al., 2023
- **VIMA** Jiang et al., 2023

## Pros

- Adaptive

# Grounding Language to Embeddings



| | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." | | |
|---|---|---|---|---|---|---|---|---|
| **Target Representation** | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills | Natural Language | Vision | Vector |
| | in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you | Steps: Step 1: chop tree block Step 2: craft wooden stick... | Images | Model Activations |

**Discrete More Structure More Bias** ←————————————————————————————————————→ **Continuous Less Structure More Variance**

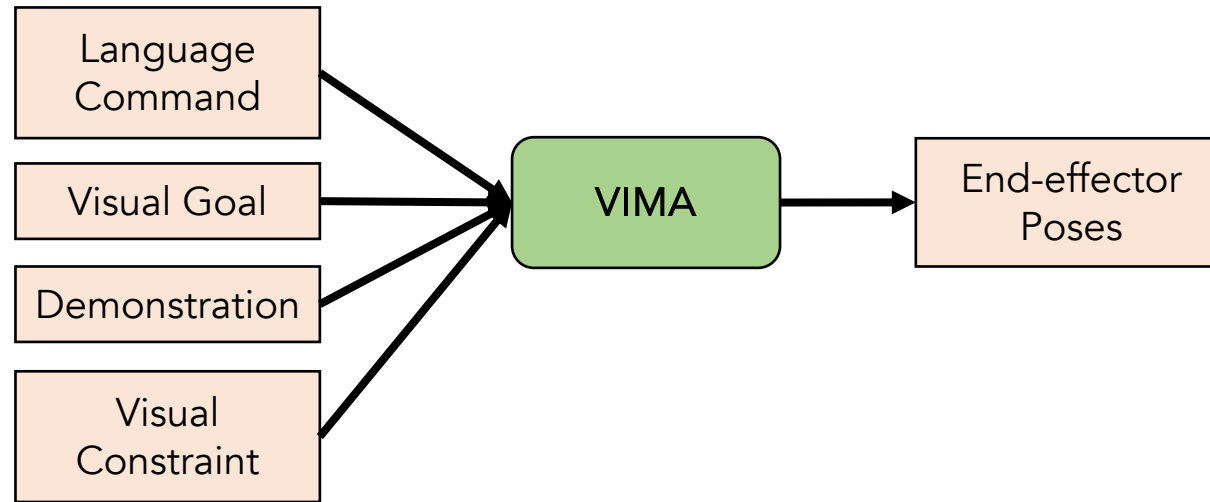| LEFT Hsu et al., 2023 | Lang2LTL Liu et al., 2023b | LLM+P Liu et al., 2023 | Code as Policies Liang et al., 2022 | SayCan Ahn et al. 2022 | VLP Du et al., 2023 | VIMA Jiang et al., 2023 |
|---|---|---|---|---|---|---|

## Pros

- Adaptive

## Cons

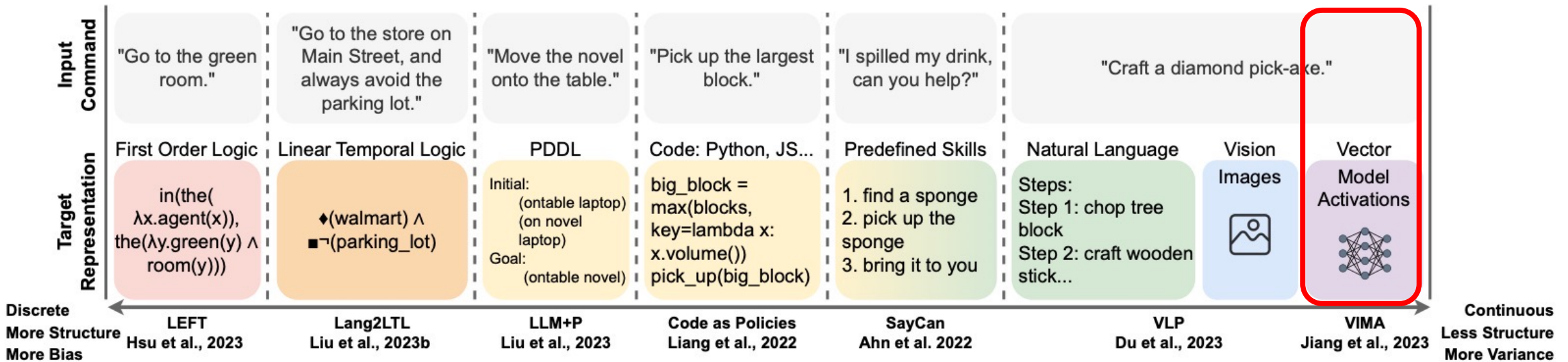- Large training set and compute
- Possibly incorrect actions

# Grounding Language to Embeddings: VIMA



VIMA

- Tokenize multimodal input

- Transformer architecture

- Output end-effector poses

# Grounding Language to Embeddings



| | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills | Natural Language | Vision | Vector |
|---|---|---|---|---|---|---|---|---|

**Input Command:**
"Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." | | |

**Target Representation:**

First Order Logic:
in(the($\lambda$x.agent(x)), the($\lambda$y.green(y) $\wedge$ room(y)))

Linear Temporal Logic:
$\blacklozenge$(walmart) $\wedge$ $\blacksquare$¬(parking_lot)

PDDL:
Initial:
(ontable laptop)
(on novel laptop)
Goal:
(ontable novel)

Code:
big_block = max(blocks, key=lambda x: x.volume())
pick_up(big_block)

Predefined Skills:
1. find a sponge
2. pick up the sponge
3. bring it to you

Natural Language:
Steps:
Step 1: chop tree block
Step 2: craft wooden stick...

Vision: Images

Vector: Model Activations

**Discrete / More Structure / More Bias** ← → **Continuous / Less Structure / More Variance**

- LEFT — Hsu et al., 2023
- Lang2LTL — Liu et al., 2023b
- LLM+P — Liu et al., 2023
- Code as Policies — Liang et al., 2022
- SayCan — Ahn et al. 2022
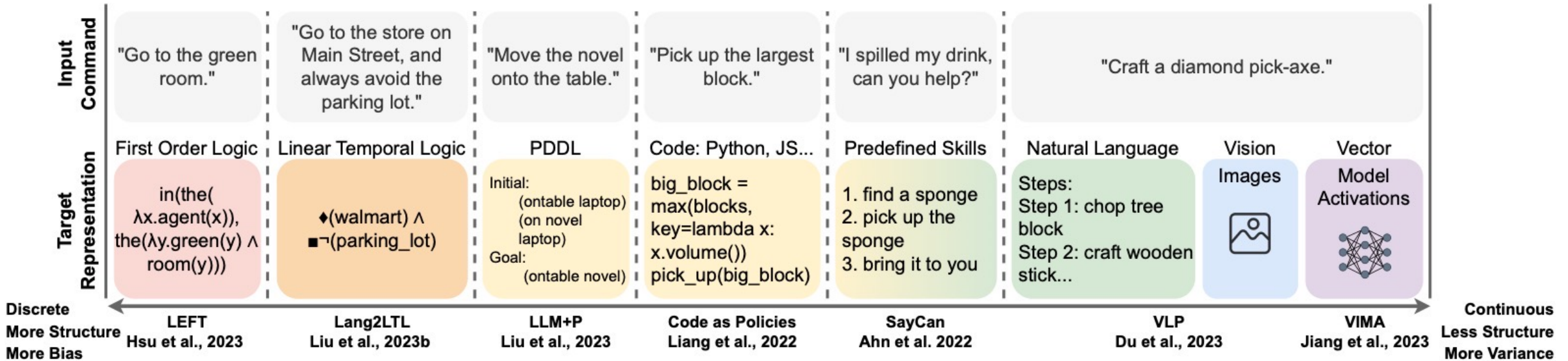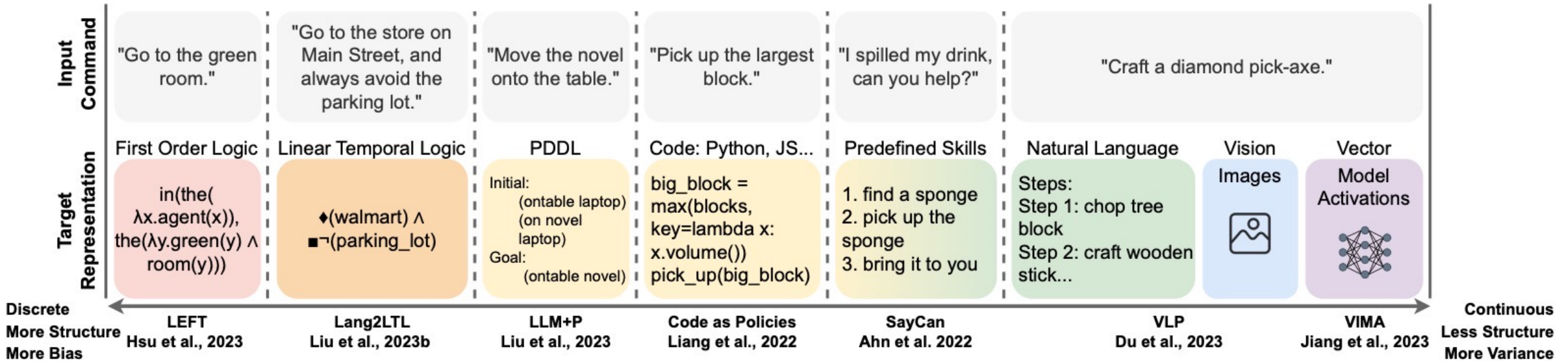- VLP — Du et al., 2023
- VIMA — Jiang et al., 2023

## More Papers

- Octo: An Open-Source Generalist Robot Policy [Octo Model Team 2024]
  Open X-Embodiment: Robotic Learning Datasets and RT-X Models [Open X-Embodiment Collaboration 2024]
- RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control [Brohan et al. 2023]
- RT-1: Robotics Transformer for Real-World Control at Scale [Brohan et al. 2023]
- PaLM-E: an Embodied Multimodal Language Model [Driess et al. 2023]
- Vision-Language Foundation Models as Effective Robot Imitators [Li et al. 2023]
- GATO: A Generalist Agent [Reed et al. 2022]
- Perceiver-Actor: A Multi-Task Transformer for Robotic Manipulation [Shridhar et al. 2022]
- Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos [Baker et al. 2022]
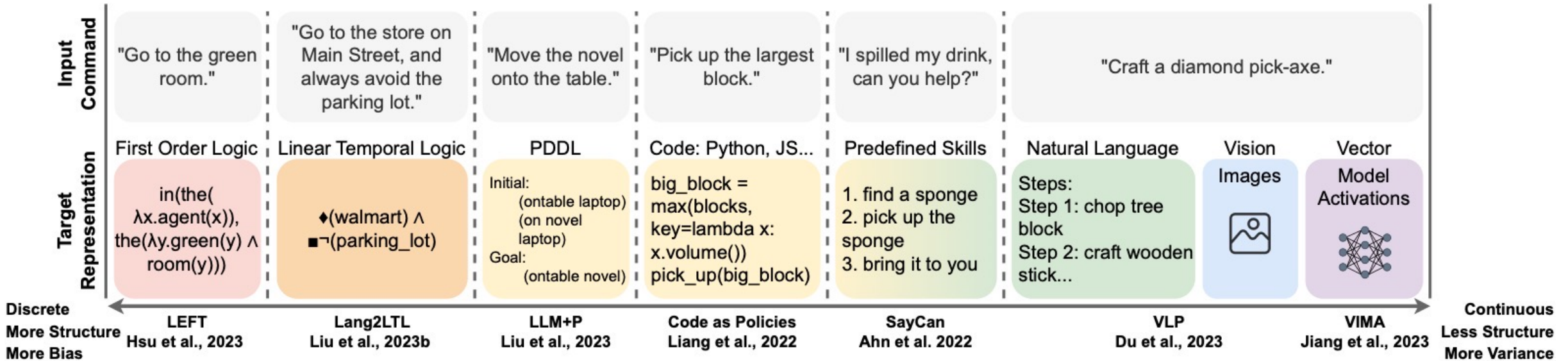
# Language Grounding for Robots



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." | |
| **Target Representation** | First Order Logic<br><br>in(the(<br>λx.agent(x)),<br>the(λy.green(y) ∧<br>room(y))) | Linear Temporal Logic<br><br>♦(walmart) ∧<br>■¬(parking_lot) | PDDL<br><br>Initial:<br>(ontable laptop)<br>(on novel<br>laptop)<br>Goal:<br>(ontable novel) | Code: Python, JS...<br><br>big_block =<br>max(blocks,<br>key=lambda x:<br>x.volume())<br>pick_up(big_block) | Predefined Skills<br><br>1. find a sponge<br>2. pick up the sponge<br>3. bring it to you | Natural Language<br><br>Steps:<br>Step 1: chop tree block<br>Step 2: craft wooden stick... | Vision<br>Images / Vector<br>Model Activations |

**Discrete**
**More Structure**
**More Bias**

**LEFT**
Hsu et al., 2023

**Lang2LTL**
Liu et al., 2023b

**LLM+P**
Liu et al., 2023

**Code as Policies**
Liang et al., 2022

**SayCan**
Ahn et al. 2022

**VLP**
Du et al., 2023

**VIMA**
Jiang et al., 2023

**Continuous**
**Less Structure**
**More Variance**

# Language Grounding for Robots



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." | |
| **Target Representation** | First Order Logic | Linear Temporal Logic | PDDL | Code: Python, JS... | Predefined Skills | Natural Language | Vision | Vector |
| | in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | ♦(walmart) ∧ ■¬(parking_lot) | Initial: (ontable laptop) (on novel laptop) Goal: (ontable novel) | big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | 1. find a sponge 2. pick up the sponge 3. bring it to you | Steps: Step 1: chop tree block Step 2: craft wooden stick... | Images | Model Activations |

**Discrete** — **Continuous**
**More Structure** — **Less Structure**
**More Bias** — **More Variance**

| LEFT Hsu et al., 2023 | Lang2LTL Liu et al., 2023b | LLM+P Liu et al., 2023 | Code as Policies Liang et al., 2022 | SayCan Ahn et al. 2022 | VLP Du et al., 2023 | VIMA Jiang et al., 2023 |
|---|---|---|---|---|---|---|

## Discrete Symbols

- Logic

- Planning domain definition language (PDDL)

- Code

- Descriptions of predefined skills

# Language Grounding for Robots

| | | | | | | |
|---|---|---|---|---|---|---|
| **Input Command** | "Go to the green room." | "Go to the store on Main Street, and always avoid the parking lot." | "Move the novel onto the table." | "Pick up the largest block." | "I spilled my drink, can you help?" | "Craft a diamond pick-axe." |
| **Target Representation** | First Order Logic<br><br>in(the(λx.agent(x)), the(λy.green(y) ∧ room(y))) | Linear Temporal Logic<br><br>♦(walmart) ∧ ■¬(parking_lot) | PDDL<br><br>Initial:<br>(ontable laptop)<br>(on novel laptop)<br>Goal:<br>(ontable novel) | Code: Python, JS...<br><br>big_block = max(blocks, key=lambda x: x.volume()) pick_up(big_block) | Predefined Skills<br><br>1. find a sponge<br>2. pick up the sponge<br>3. bring it to you | Natural Language    Vision    Vector<br><br>Steps:<br>Step 1: chop tree block<br>Step 2: craft wooden stick...   Images   Model Activations |

**Discrete** → **Continuous**
More Structure / More Bias → Less Structure / More Variance

| LEFT<br>Hsu et al., 2023 | Lang2LTL<br>Liu et al., 2023b | LLM+P<br>Liu et al., 2023 | Code as Policies<br>Liang et al., 2022 | SayCan<br>Ahn et al. 2022 | VLP<br>Du et al., 2023 | VIMA<br>Jiang et al., 2023 |

**Discrete Symbols**

- Logic
- Planning domain definition language (PDDL)
- Code
- Descriptions of predefined skills

**High-dimensional Embeddings**

- Language and image subgoals
- Neural embeddings

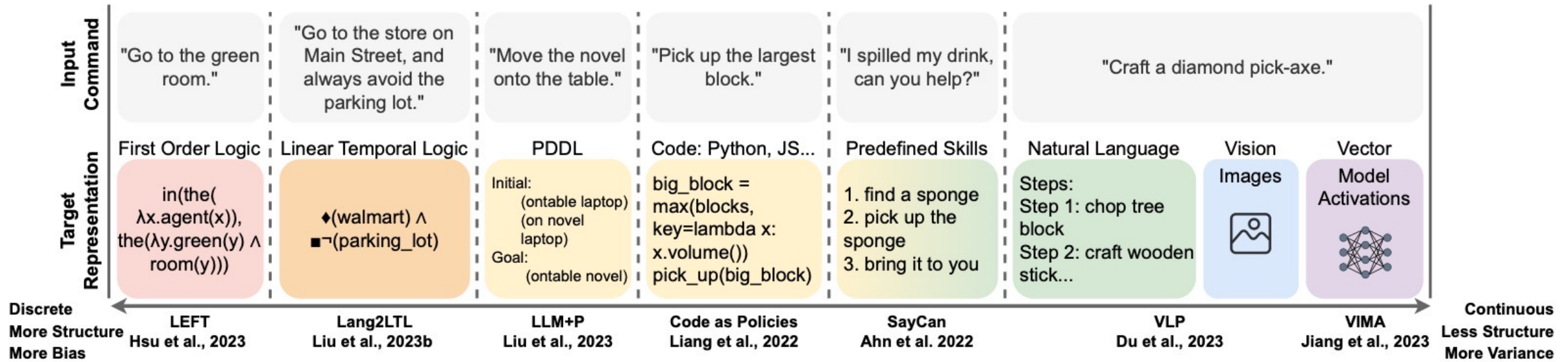# Open Problems and Future Directions

# Open Problems and Future Directions

- **Neuro-symbolic Approach**

    - POMDP and PDDL planners

    - Deep learning models with generalizable representations

    - E.g., Jointly learn symbols in the embedding space and skills

# Open Problems and Future Directions

- **Neuro-symbolic Approach**

  - POMDP and PDDL planners

  - Deep learning models with generalizable representations

  - E.g., Jointly learn symbols in the embedding space and skills

- **Multimodal Dataset**

  - E.g., text, audio, RGB images, point clouds, voxels, videos, demonstrations

  - Semantically diverse

# Open Problems and Future Directions

- **Neuro-symbolic Approach**

  - POMDP and PDDL planners

  - Deep learning models with generalizable representations

  - E.g., Jointly learn symbols in the embedding space and skills

- **Multimodal Dataset**

  - E.g., text, audio, RGB images, point clouds, voxels, videos, demonstrations

  - Semantically diverse

- **Modular Approach**

  - Existing robot modules

  - E.g., SLAM, motion planning and object detection

# Open Problems and Future Directions

- **Neuro-symbolic Approach**

  - POMDP and PDDL planners

  - Deep learning models with generalizable representations

  - E.g., Jointly learn symbols in the embedding space and skills

- **Multimodal Dataset**

  - E.g., text, audio, RGB images, point clouds, voxels, videos, demonstrations

  - Semantically diverse

- **Modular Approach**

  - Existing robot modules

  - E.g., SLAM, motion planning and object detection

- **Verification and Safety**

  - Formal methods

# Conclusion



A Survey of Robotic Language Grounding: Tradeoffs between Symbols and Embeddings

https://arxiv.org/abs/2405.13245

Jason Xinyu Liu

https://jasonxyliu.github.io/

xinyu_liu@brown.edu

Poster Location:
E15

45