



Lang2LTL-2: Grounding Spatiotemporal Navigation Commands Using Large Language and Vision-Language Models

Jason Xinyu Liu



<https://jasonxyliu.github.io>

AAAI 2024 Fall Symposium UR-RAD

November 8st, 2024



<https://spatiotemporal-ground.github.io>

Grounding Spatiotemporal Language

go to the white car near the dumpster
exactly three times, in addition avoid stairs
in front of the apartment



walk to the chair in front of the bookshelf
but only after the kitchen counter



Grounding Language to Structure

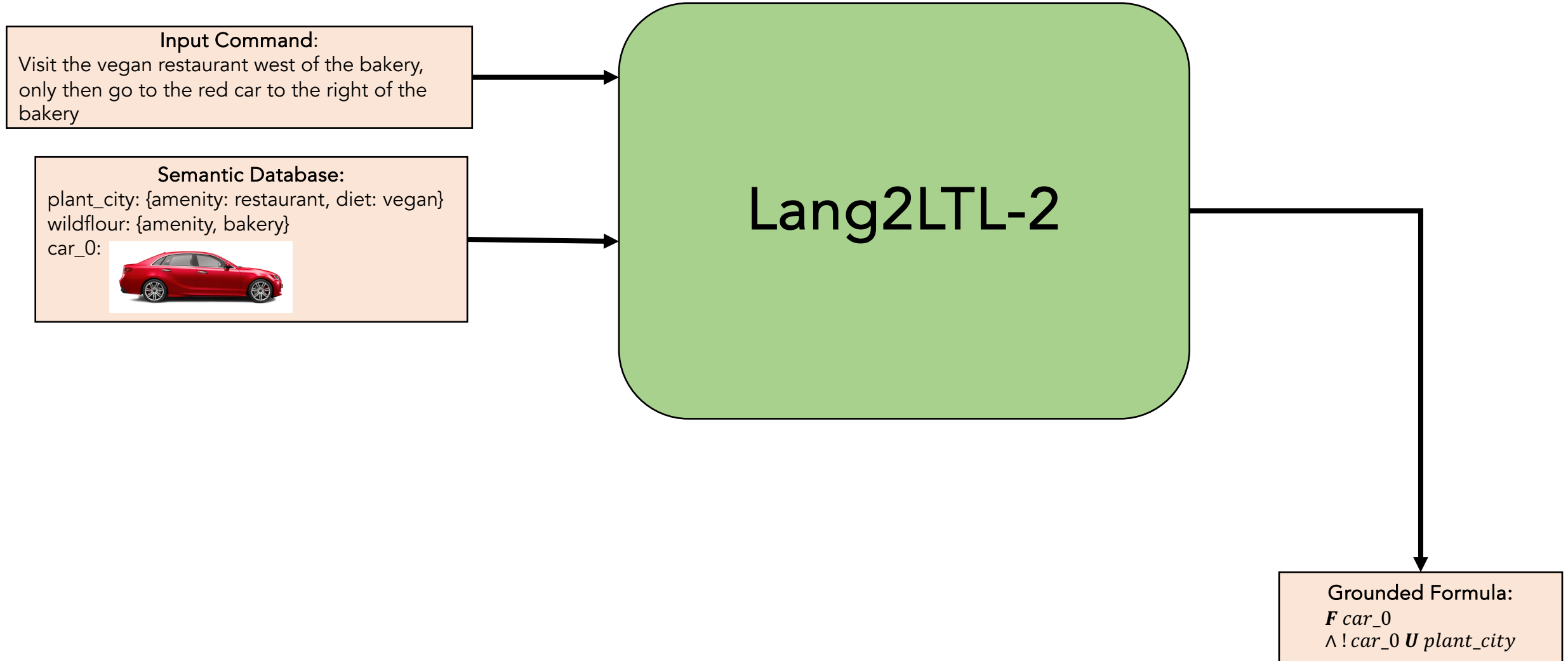
Input

- Spatiotemporal navigation command
- Multimodal semantic map: text + images

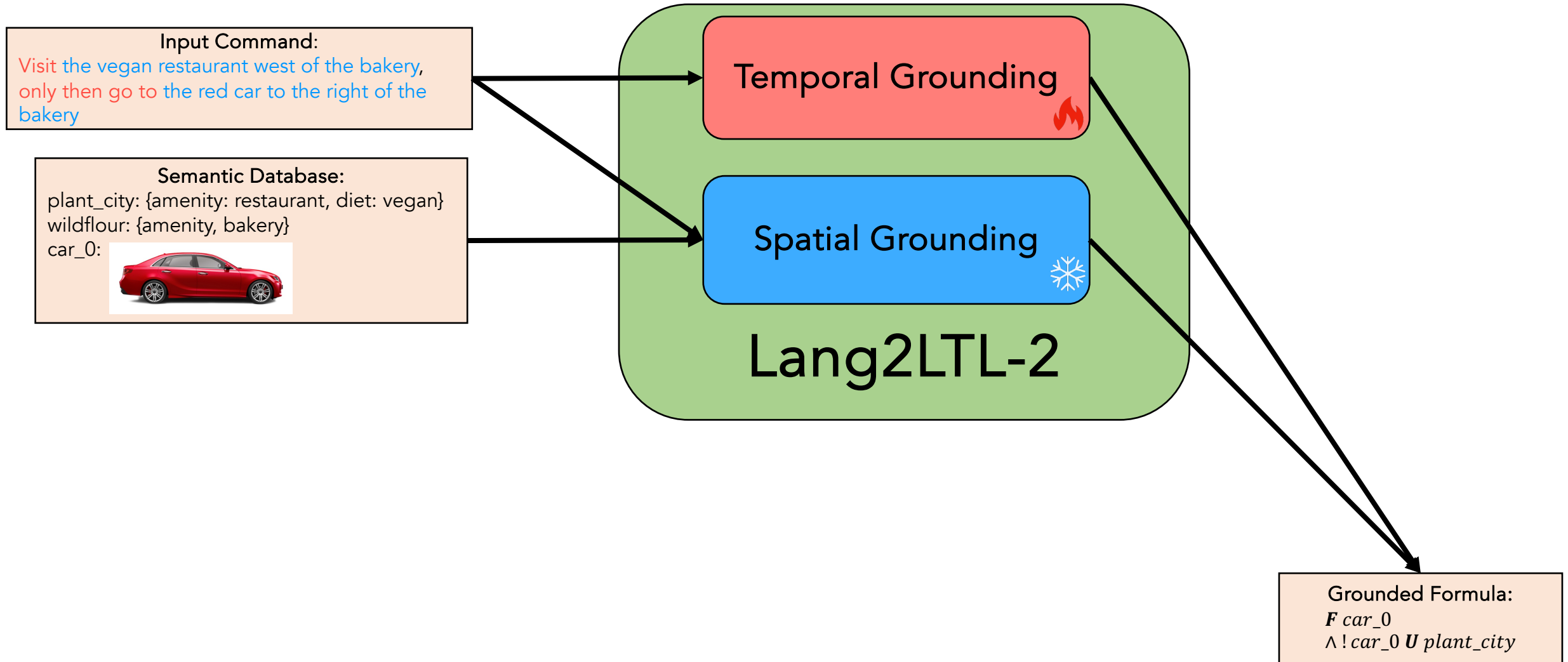
Output

- LTL formula whose propositions are grounded to real-world landmarks

Lang2LTL-2: A Modular Grounding System



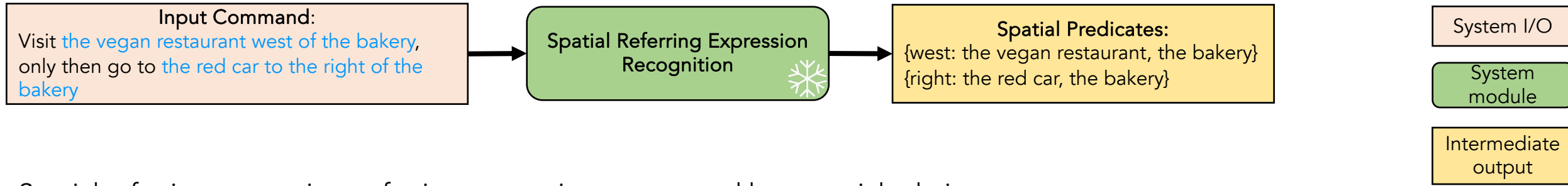
Lang2LTL-2: A Modular Grounding System



Lang2LTL-2: A Modular Grounding System

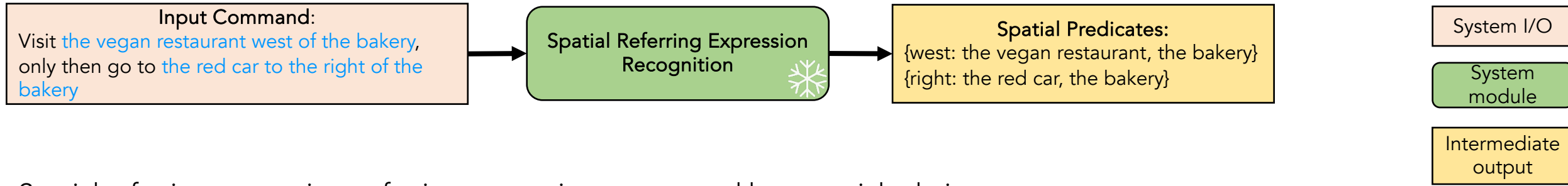


Lang2LTL-2: A Modular Grounding System



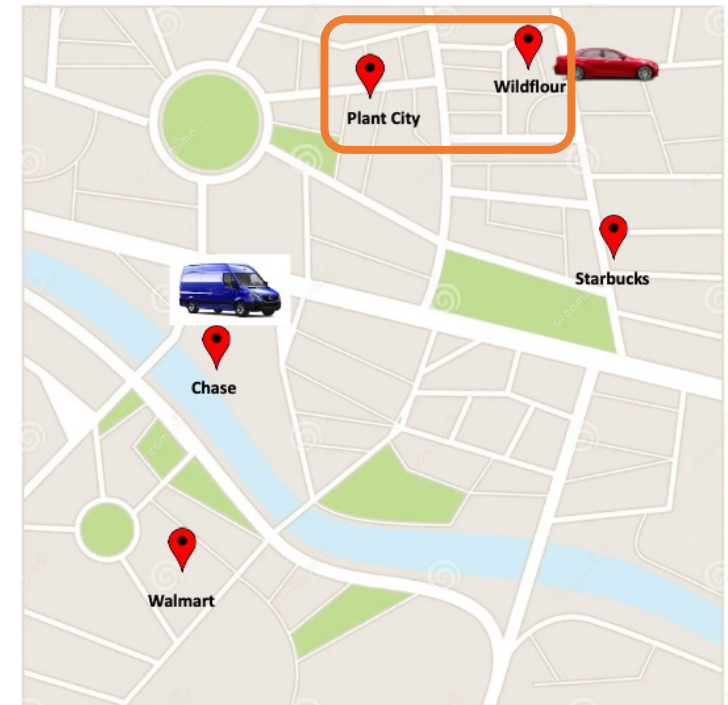
Spatial referring expression: referring expressions connected by a spatial relation

Lang2LTL-2: A Modular Grounding System

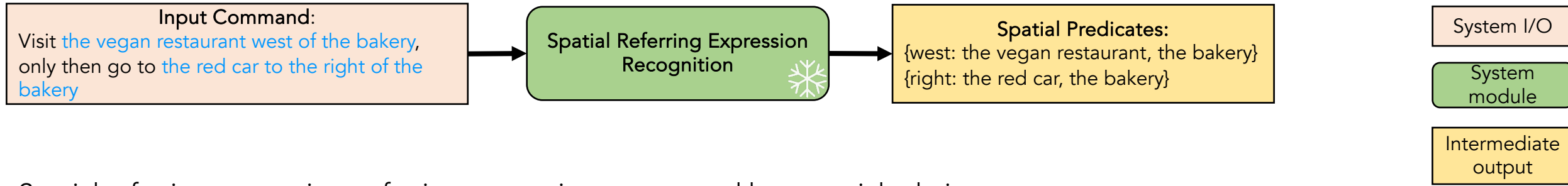


Spatial referring expression: referring expressions connected by a spatial relation

- the vegan restaurant west of the bakery

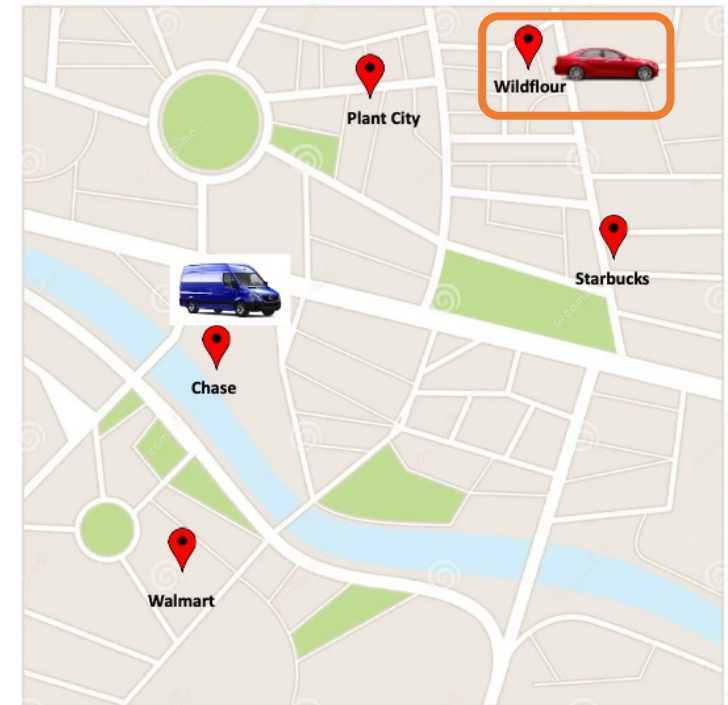


Lang2LTL-2: A Modular Grounding System

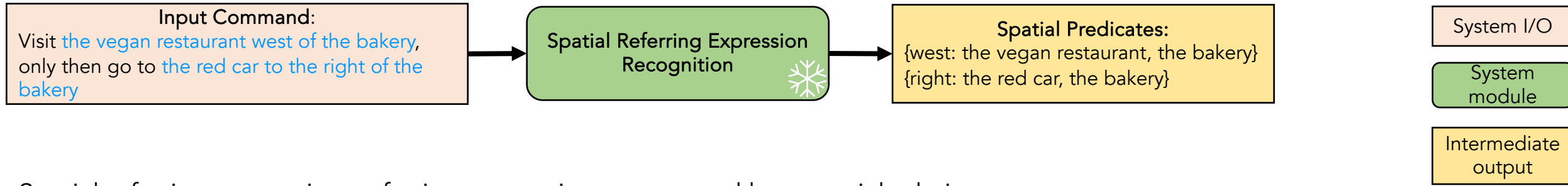


Spatial referring expression: referring expressions connected by a spatial relation

- the vegan restaurant west of the bakery
- the red car to the right of the bakery



Lang2LTL-2: A Modular Grounding System

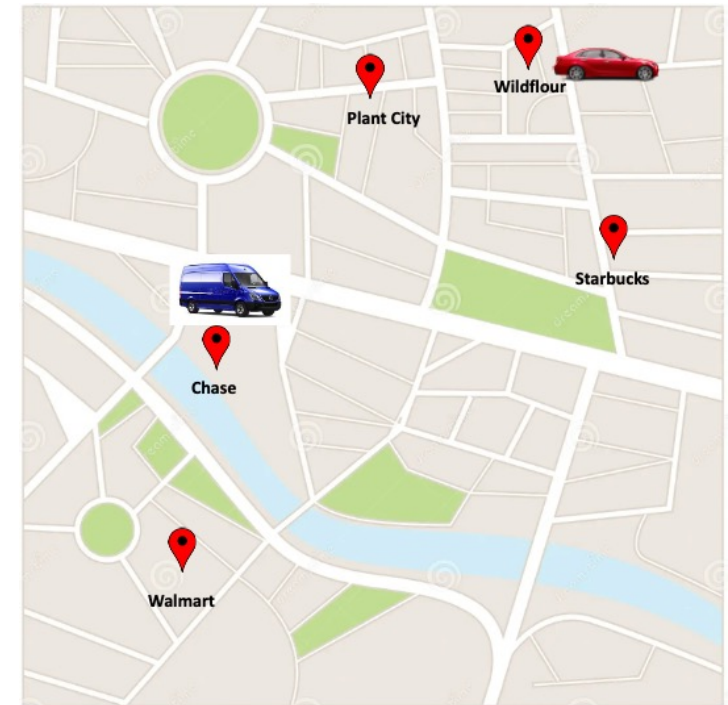


Spatial referring expression: referring expressions connected by a spatial relation

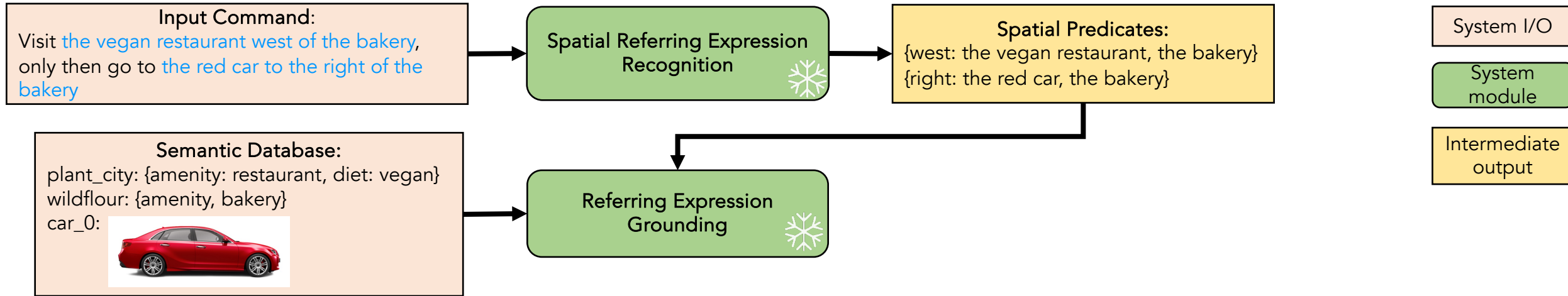
- the vegan restaurant west of the bakery
- the red car to the right of the bakery

Spatial referring expression recognition module

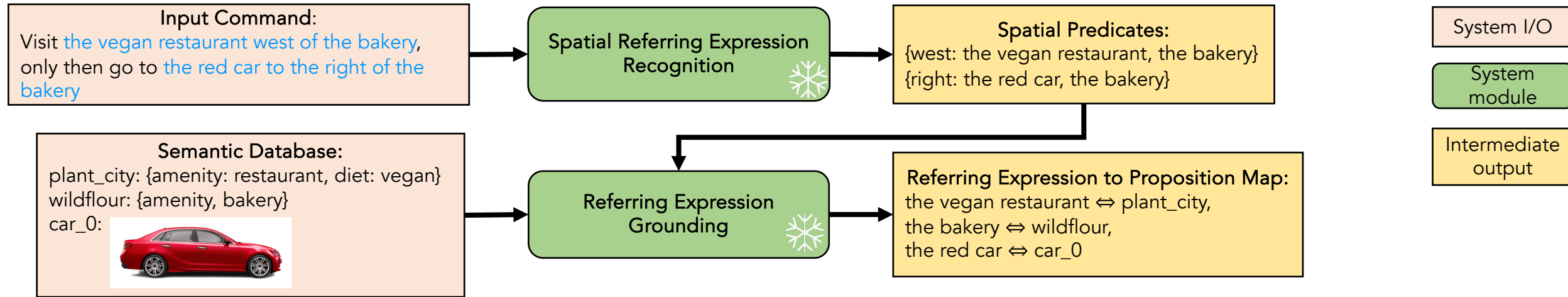
- In-context learning with GPT-4



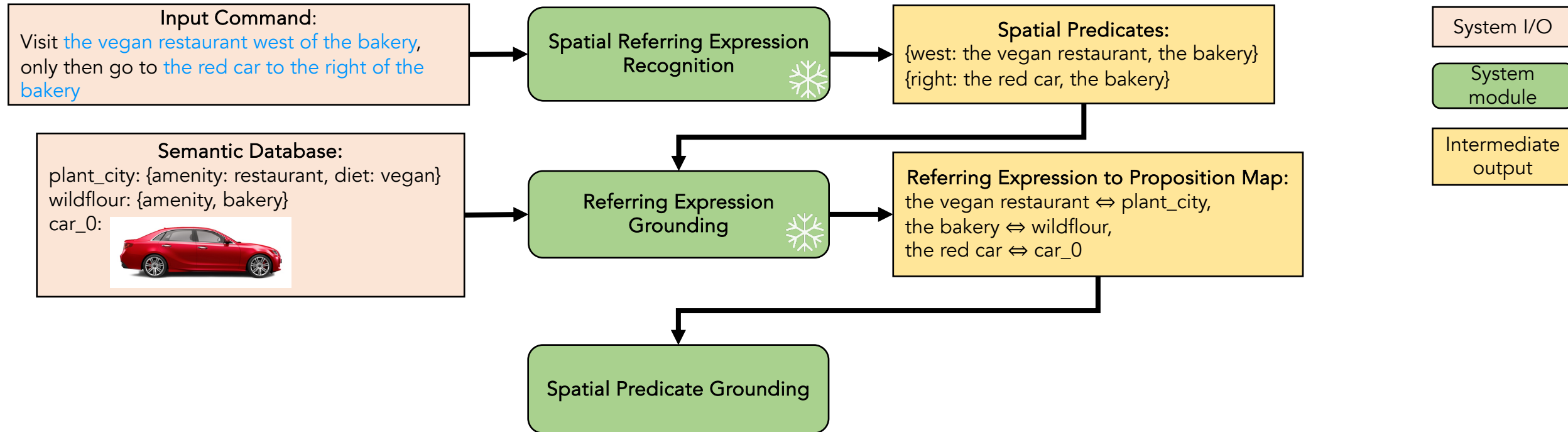
Lang2LTL-2: A Modular Grounding System



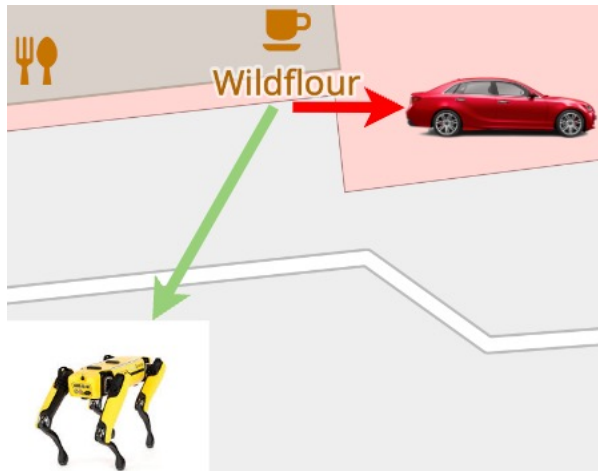
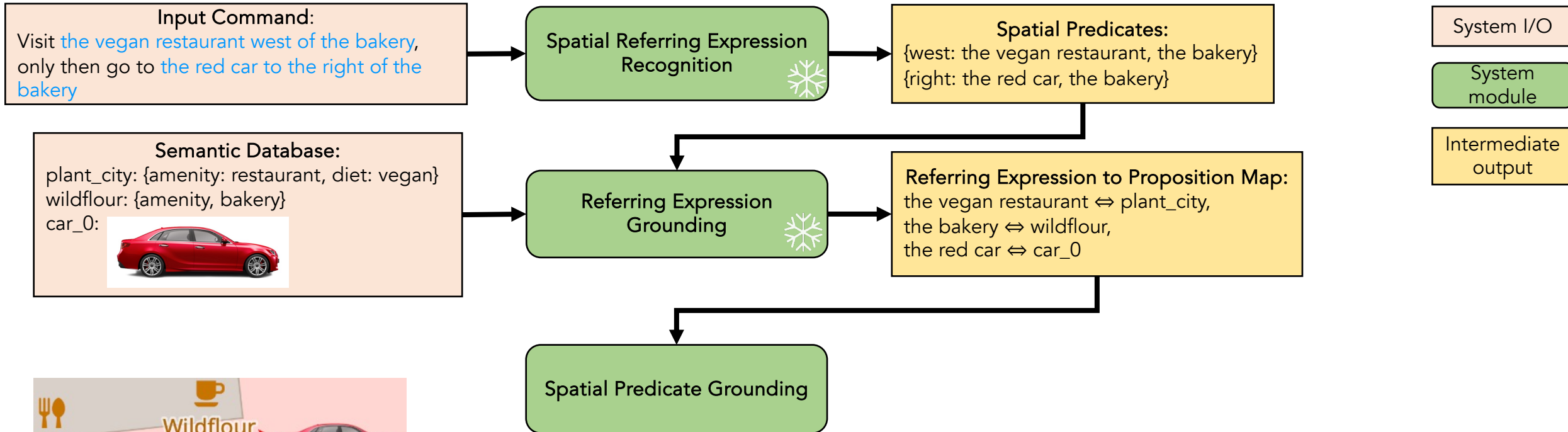
Lang2LTL-2: A Modular Grounding System



Lang2LTL-2: A Modular Grounding System

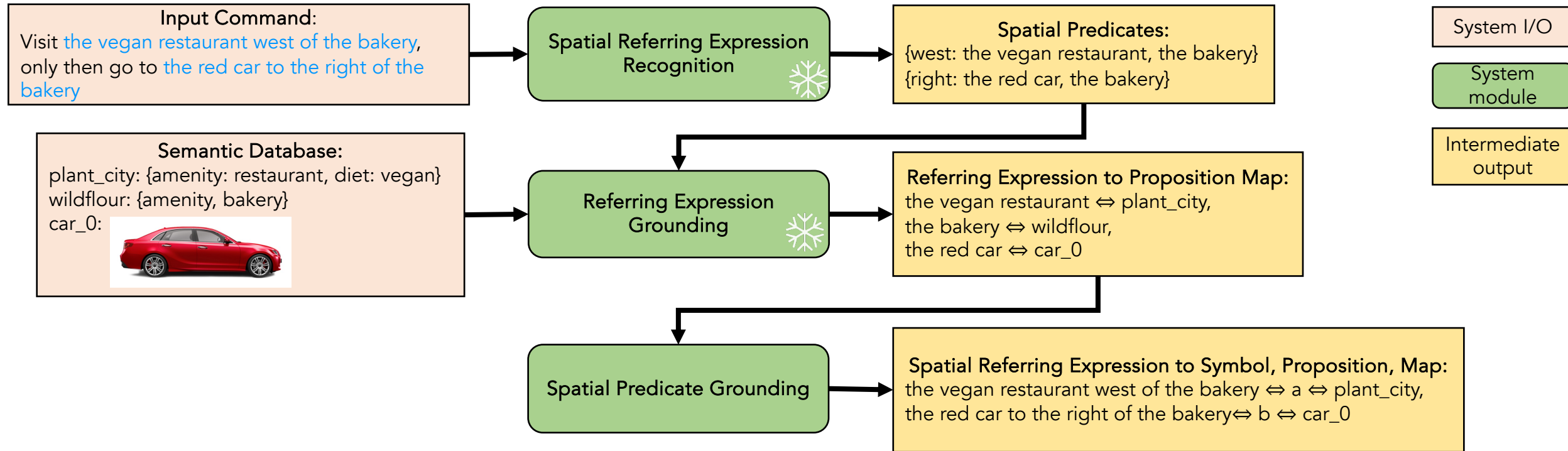


Lang2LTL-2: A Modular Grounding System

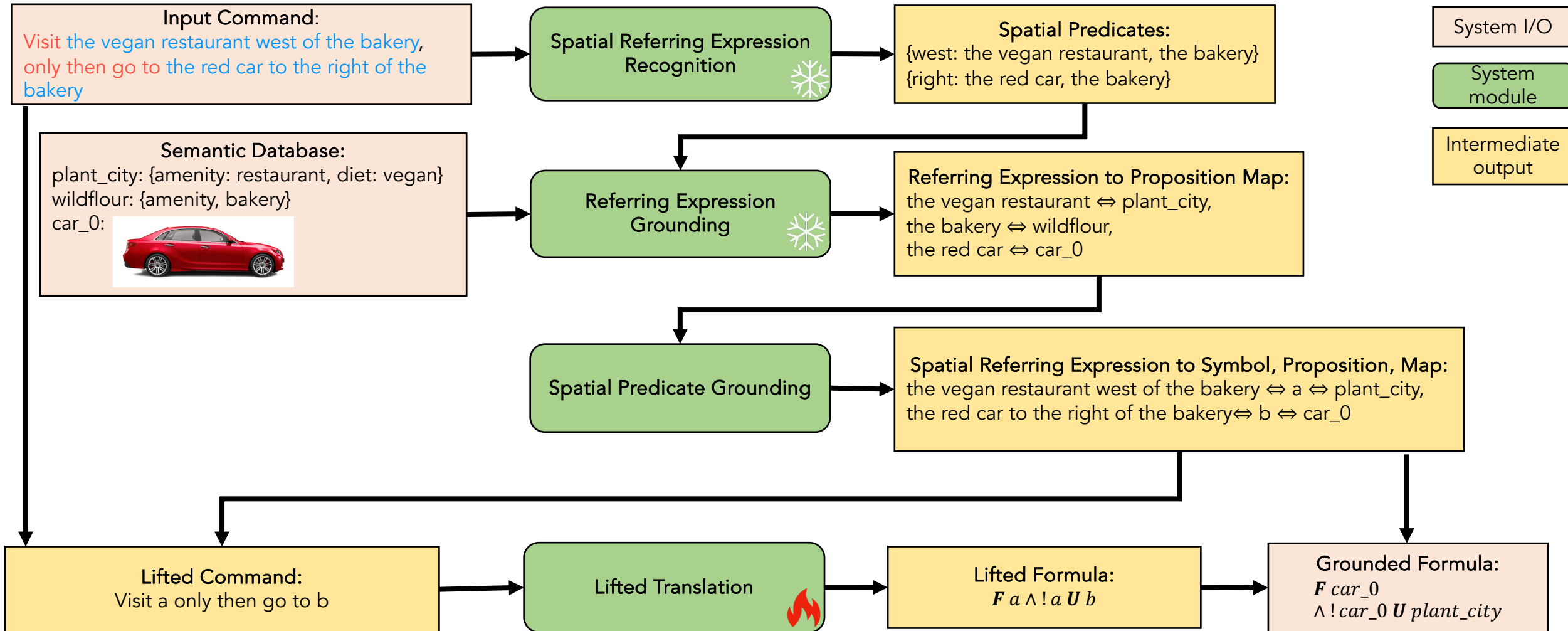


- Spatial referring expression: the red car to the right of the bakery
- Spatial predicate: {right: the red car, the bakery}
- Spatial relation: to the right of

Lang2LTL-2: A Modular Grounding System



Lang2LTL-2: A Modular Grounding System



Lang2LTL-2: Modular Evaluation

Module	Accuracy				
	<i>City 1</i> (9 landmarks)	<i>City 2</i> (34 landmarks)	<i>City 3</i> (44 landmarks)	<i>City 4</i> (175 landmarks)	<i>Average</i>
SRER					
REG					
SPG					
LT					

Lang2LTL-2: Modular Evaluation

Module	Accuracy				
	<i>City 1</i> (9 landmarks)	<i>City 2</i> (34 landmarks)	<i>City 3</i> (44 landmarks)	<i>City 4</i> (175 landmarks)	<i>Average</i>
SRER	99.45 ± 0.12%	99.43 ± 0.26%	99.56 ± 0.63%	99.39 ± 0.21%	99.46 ± 0.34%
REG					
SPG					
LT					

Lang2LTL-2: Modular Evaluation

Module	Accuracy					
	<i>City 1</i> (9 landmarks)	<i>City 2</i> (34 landmarks)	<i>City 3</i> (44 landmarks)	<i>City 4</i> (175 landmarks)	<i>Average</i>	
SRER	99.45 ± 0.12%	99.43 ± 0.26%	99.56 ± 0.63%	99.39 ± 0.21%	99.46 ± 0.34%	
REG	Top-1	99.68 ± 0.72%	97.98 ± 1.07%	88.74 ± 2.14%	78.35 ± 1.97%	91.19 ± 8.84%
	Top-5	100.00 ± 0.00%	100.00 ± 0.00%	99.56 ± 0.24%	99.15 ± 0.34%	99.68 ± 0.41%
	Top-10	100.00 ± 0.00%	100.00 ± 0.00%	99.70 ± 0.17%	99.98 ± 0.05%	99.92 ± 0.15%
SPG						
LT						

Lang2LTL-2: Modular Evaluation

Module		Accuracy				
		<i>City 1</i> (9 landmarks)	<i>City 2</i> (34 landmarks)	<i>City 3</i> (44 landmarks)	<i>City 4</i> (175 landmarks)	<i>Average</i>
SRER		99.45 ± 0.12%	99.43 ± 0.26%	99.56 ± 0.63%	99.39 ± 0.21%	99.46 ± 0.34%
REG	Top-1	99.68 ± 0.72%	97.98 ± 1.07%	88.74 ± 2.14%	78.35 ± 1.97%	91.19 ± 8.84%
	Top-5	100.00 ± 0.00%	100.00 ± 0.00%	99.56 ± 0.24%	99.15 ± 0.34%	99.68 ± 0.41%
	Top-10	100.00 ± 0.00%	100.00 ± 0.00%	99.70 ± 0.17%	99.98 ± 0.05%	99.92 ± 0.15%
SPG		100.00 ± 0.00%	100.00 ± 0.00%	99.53 ± 0.33%	99.35 ± 1.46%	99.72 ± 0.75%
LT						

Lang2LTL-2: Modular Evaluation

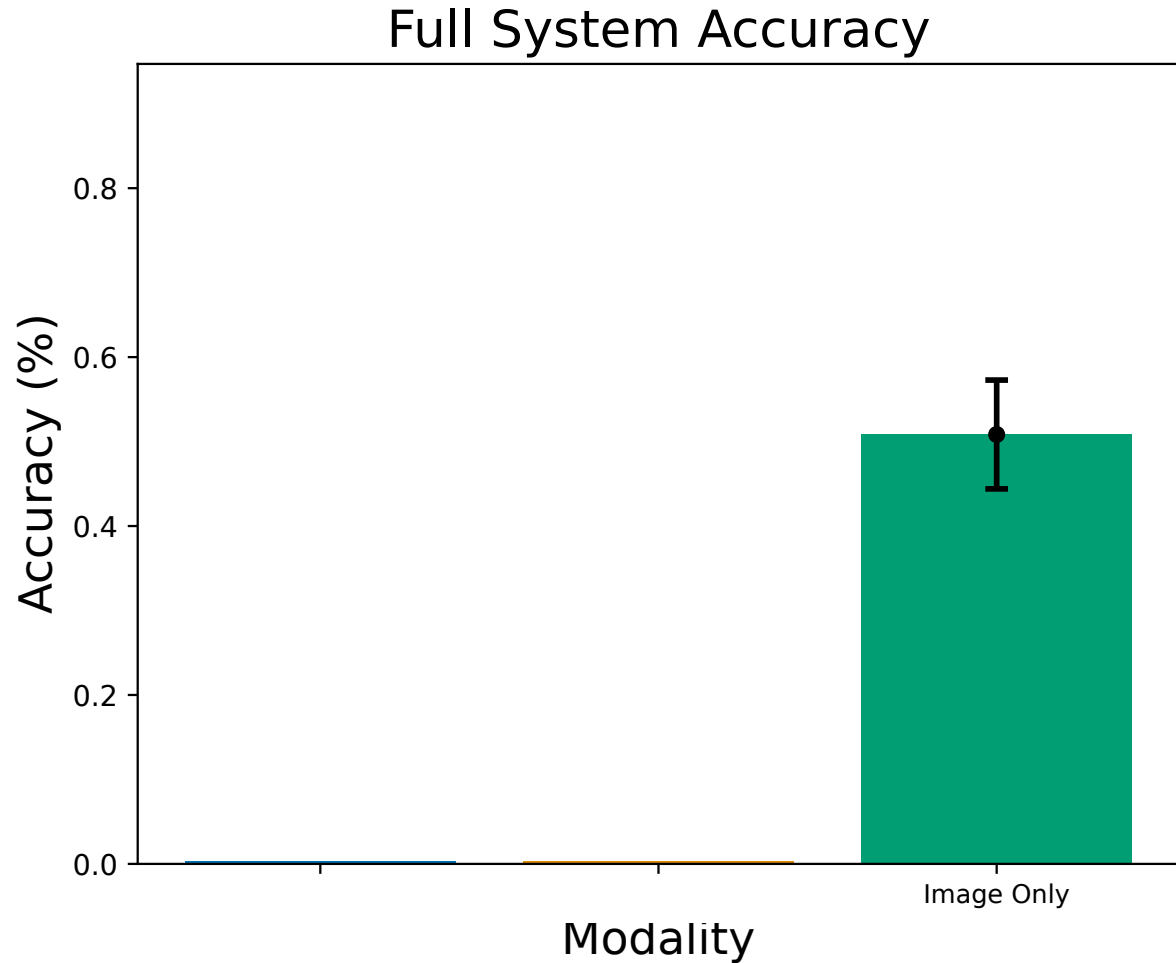
Module		Accuracy				
		City 1 (9 landmarks)	City 2 (34 landmarks)	City 3 (44 landmarks)	City 4 (175 landmarks)	Average
SRER		99.45 ± 0.12%	99.43 ± 0.26%	99.56 ± 0.63%	99.39 ± 0.21%	99.46 ± 0.34%
REG	Top-1	99.68 ± 0.72%	97.98 ± 1.07%	88.74 ± 2.14%	78.35 ± 1.97%	91.19 ± 8.84%
	Top-5	100.00 ± 0.00%	100.00 ± 0.00%	99.56 ± 0.24%	99.15 ± 0.34%	99.68 ± 0.41%
	Top-10	100.00 ± 0.00%	100.00 ± 0.00%	99.70 ± 0.17%	99.98 ± 0.05%	99.92 ± 0.15%
SPG		100.00 ± 0.00%	100.00 ± 0.00%	99.53 ± 0.33%	99.35 ± 1.46%	99.72 ± 0.75%
LT	Finetuned T5-base	99.45 ± 0.00%	99.45 ± 0.00%	99.45 ± 0.00%	99.45 ± 0.00%	99.45 ± 0.00%
	RAG-10	69.33 ± 0.25%	70.34 ± 0.13%	69.65 ± 0.58%	70.39 ± 0.84%	69.93 ± 0.62%
	RAG-50	83.79 ± 0.06%	83.93 ± 0.12%	83.75 ± 0.52%	83.93 ± 0.65%	83.85 ± 0.33%
	RAG-100	88.20 ± 0.58%	88.25 ± 1.04%	87.79 ± 0.39%	87.70 ± 0.13%	87.98 ± 0.54%

Lang2LTL-2: Full System Evaluation

- Average success rates
- 21,780 semantically diverse commands
 - 47 temporal patterns
 - 19 spatial relations
- 4 cities
- 5 seeds

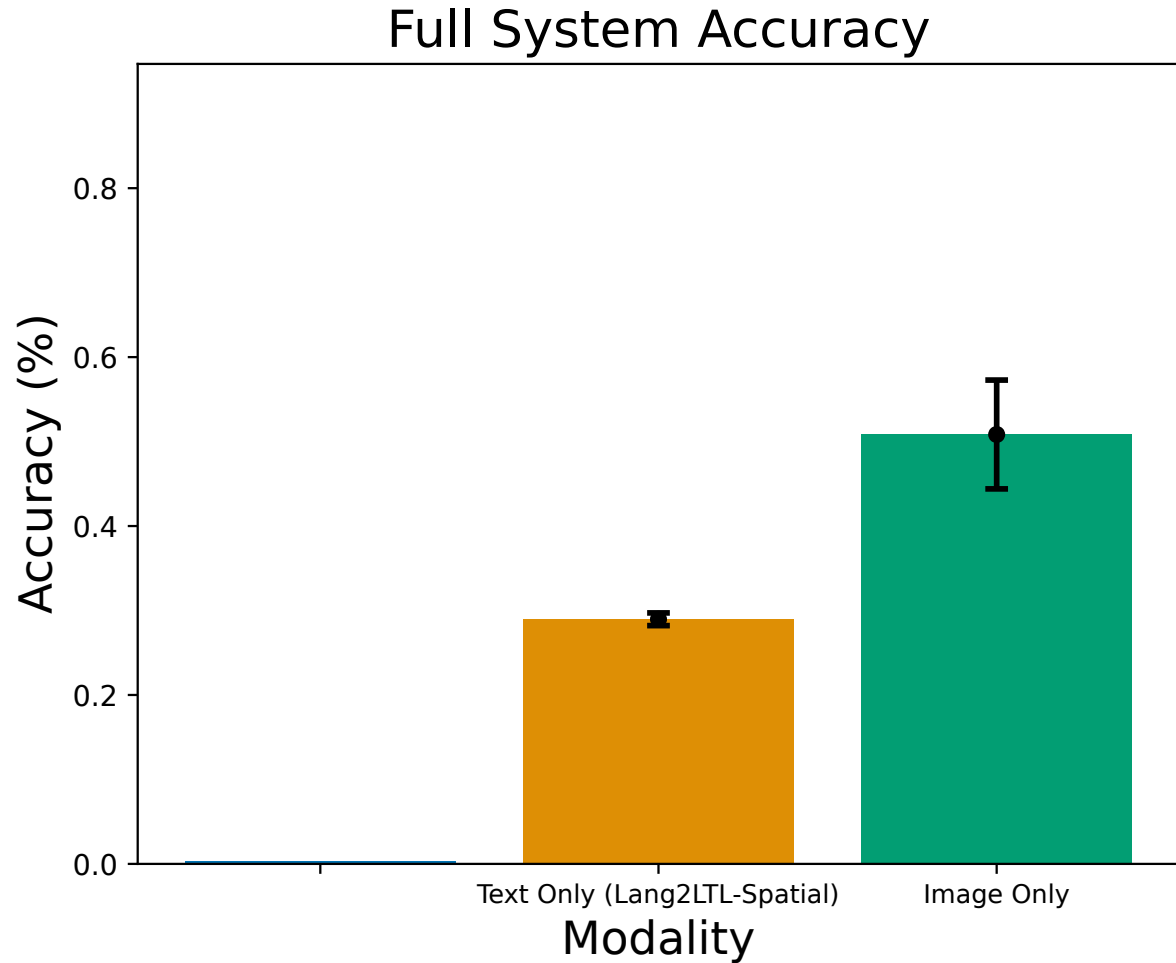
Lang2LTL-2: Full System Evaluation

- Average success rates
- 21,780 semantically diverse commands
 - 47 temporal patterns
 - 19 spatial relations
- 4 cities
- 5 seeds



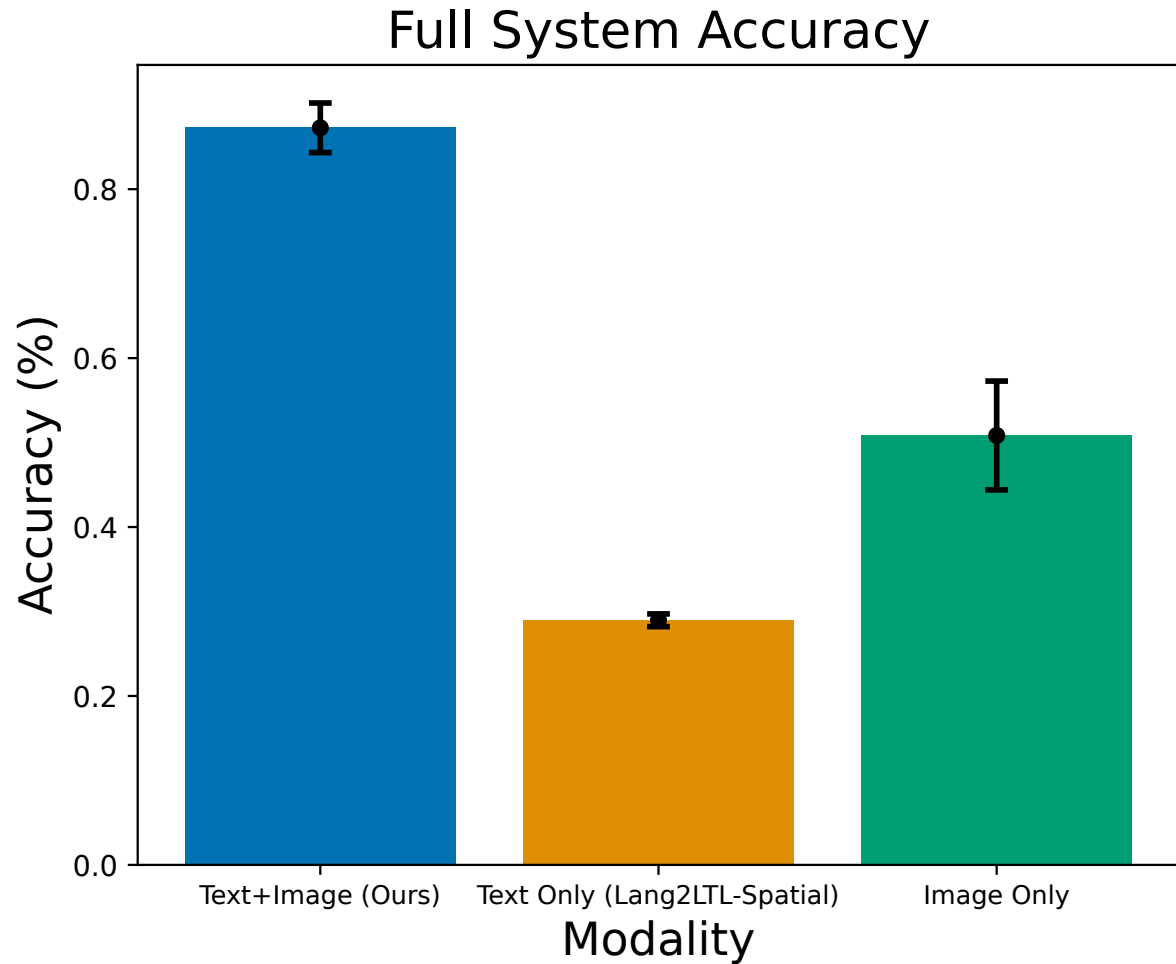
Lang2LTL-2: Full System Evaluation

- Average success rates
- 21,780 semantically diverse commands
 - 47 temporal patterns
 - 19 spatial relations
- 4 cities
- 5 seeds



Lang2LTL-2: Full System Evaluation

- Average success rates
- 21,780 semantically diverse commands
 - 47 temporal patterns
 - 19 spatial relations
- 4 cities
- 5 seeds



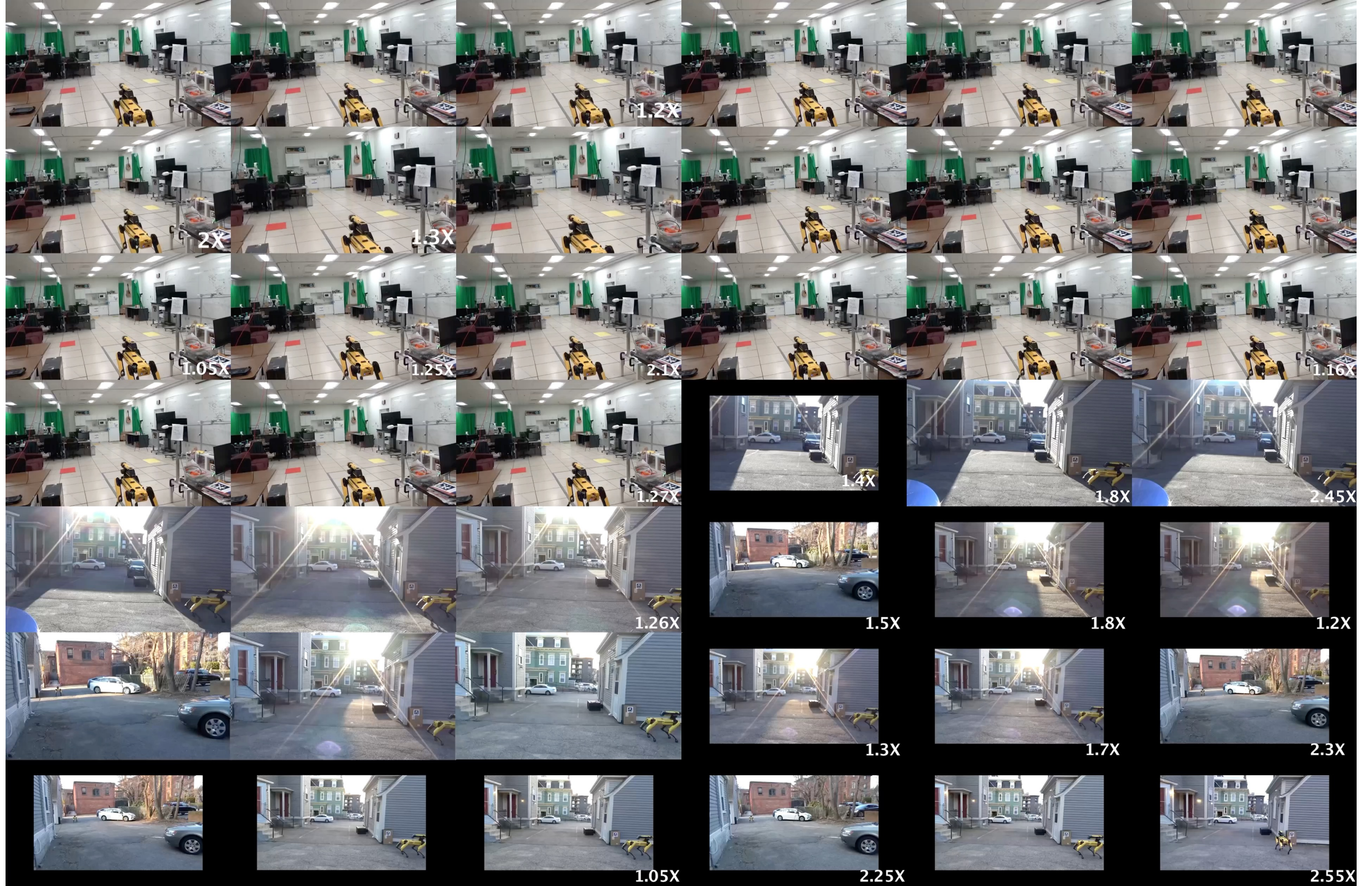
Lang2LTL-2: Robot Demonstration

go to the white car near the dumpster
exactly three times, in addition avoid stairs
in front of the apartment

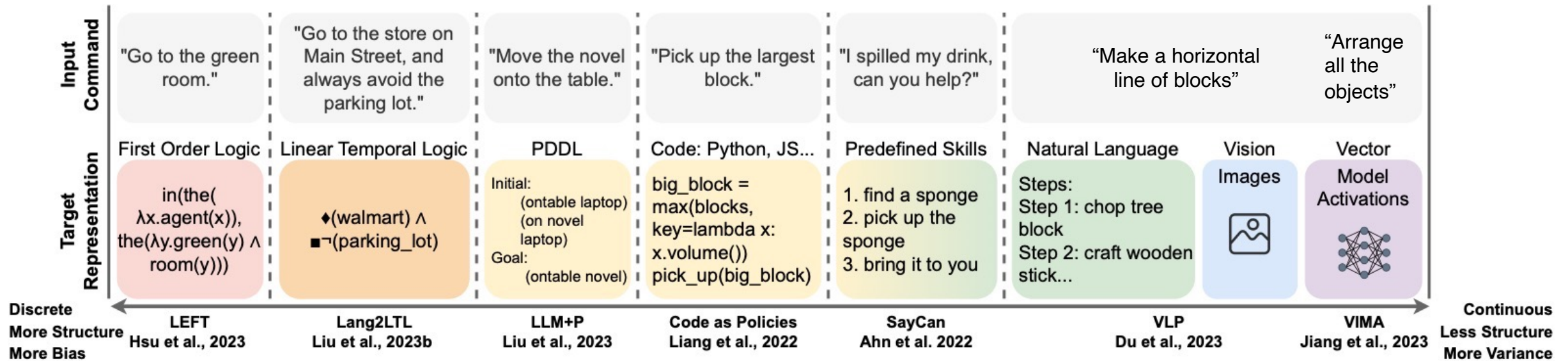


walk to the chair in front of the bookshelf
but only after the kitchen counter





Robotic Language Grounding



Symbols

- Discrete
- More Structure; More bias
- Unambiguous
- Verifiable
- Interpretable

High-dimensional Embeddings

- Continuous
- Less structure; More variance
- Adaptive



Lang2LTL-2: Grounding Spatiotemporal Language

Contributions

- Grounding in novel environments without retraining on language data: 93.53% success rate
- New benchmark of 21,780 semantically diverse commands: 47 temporal + 19 spatial
- Use multimodal semantic map
- Deployed in indoor and outdoor environments



<https://jasonxyliu.github.io>



<https://spatiotemporal-ground.github.io>



A Survey of Robotic
Language Grounding