



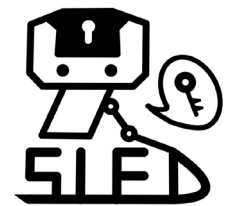
# DANLI: Deliberative Agent for Following Natural Language Instructions

Yichi Zhang, Jianing Yang, Jiayi Pan, Shane Storcks, Nikhil Devraj

Ziqiao Ma, Keunwoo Peter Yu, Yuwei Bao, and Joyce Chai

Situated Language and Embodied Dialog (SLED) Lab

University of Michigan



# Problem Definition

- TEACH<sup>[1]</sup>: The agent is given a dialog history as input, and is expected to execute a sequence of actions to achieve the goal set out by the human commander.

Please boil a potato.

Is there another pot somewhere?

You could try filling the cup with water and emptying it into the pot.

Good thinking! Thank you for that suggestion.

t=54 t=55 t=59 t=60 t=62 t=63 t=66 t=72 t=73 t=74

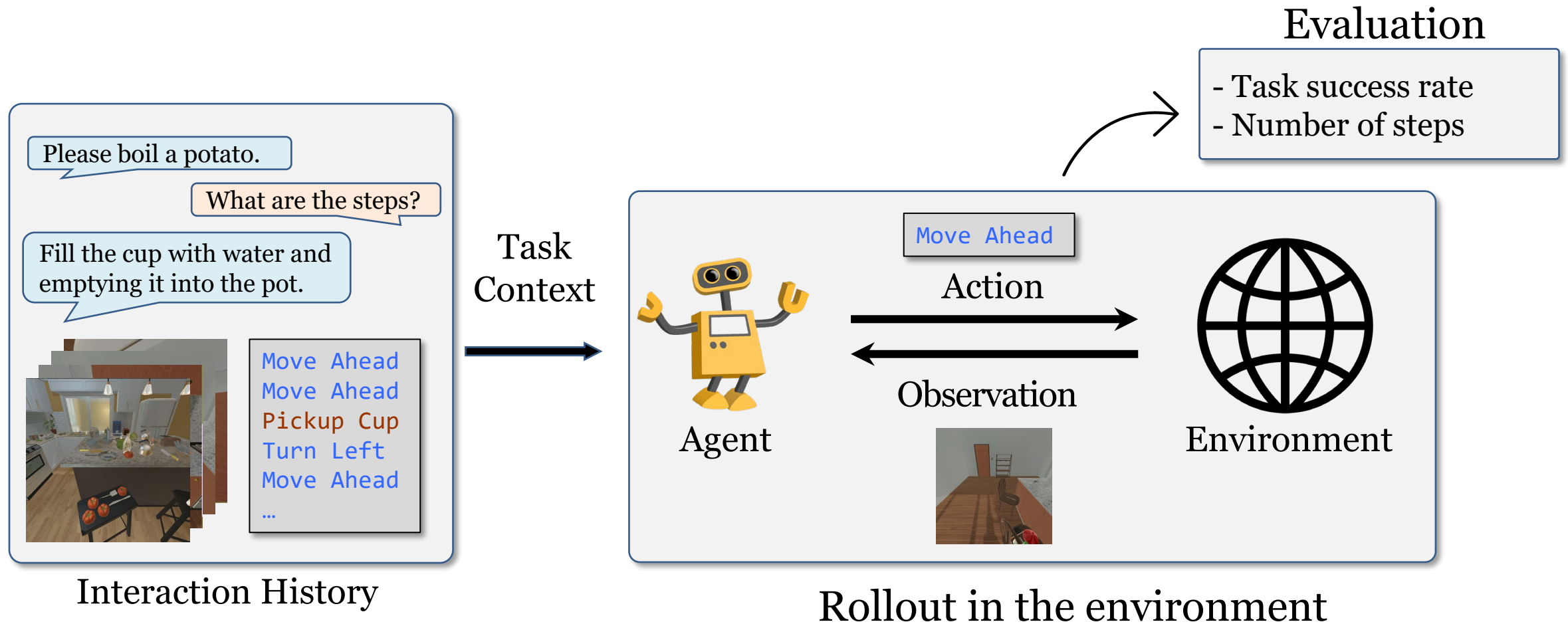
Find cup    Fill cup with water    Transfer water to pot    Boil water and add potato

## TEACH Tasks

WATER PLANT  
MAKE COFFEE  
CLEAN ALL X  
PUT ALL X ON Y  
BOIL POTATO  
MAKE PLATE OF TOAST  
N SLICES OF X IN Y  
PUT ALL X IN ONE Y  
N COOKED X SLICES IN Y  
PREPARE SANDWICH  
PREPARE SALAD  
PREPARE BREAKFAST

[1] Padmakumar A, Thomason J, Shrivastava A, Lange P, Narayan-Chen A, Gella S, PIRAMUTHU R, TUR G, HAKKANI-TUR D. Teach: Task-driven embodied agents that chat. AAAI 2022.

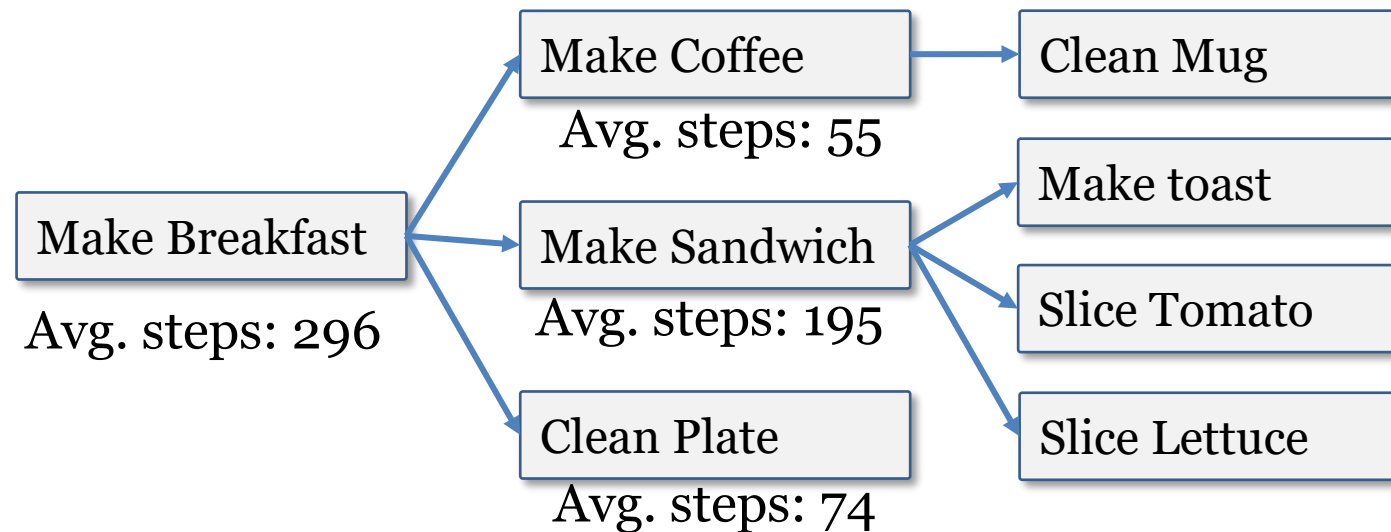
# Execution From Dialog History (EDH)



# Challenges

---

- Multi-modal situation and task understanding
- Partially observable and high-dimensional state
- Long-horizon and compositional tasks



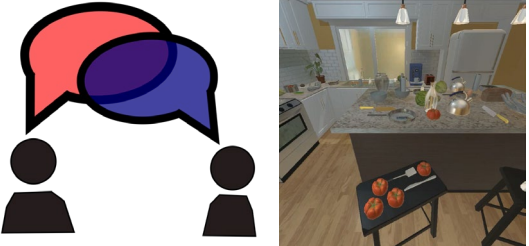
# Motivation: From Reactive to Deliberative

Pashevich et al. 2021  
Bulks et al. 2021  
Zhang and Chai. 2021  
Zhou et al. 2021  
Singh et al. 2020  
...

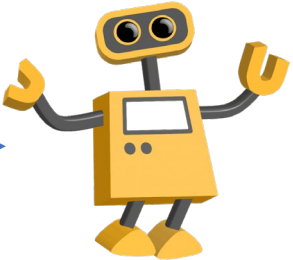
~ 7% success rate  
on TEACH-EDH 🥲



*Reactive*



Situation



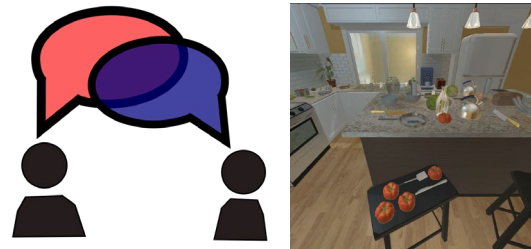
Action



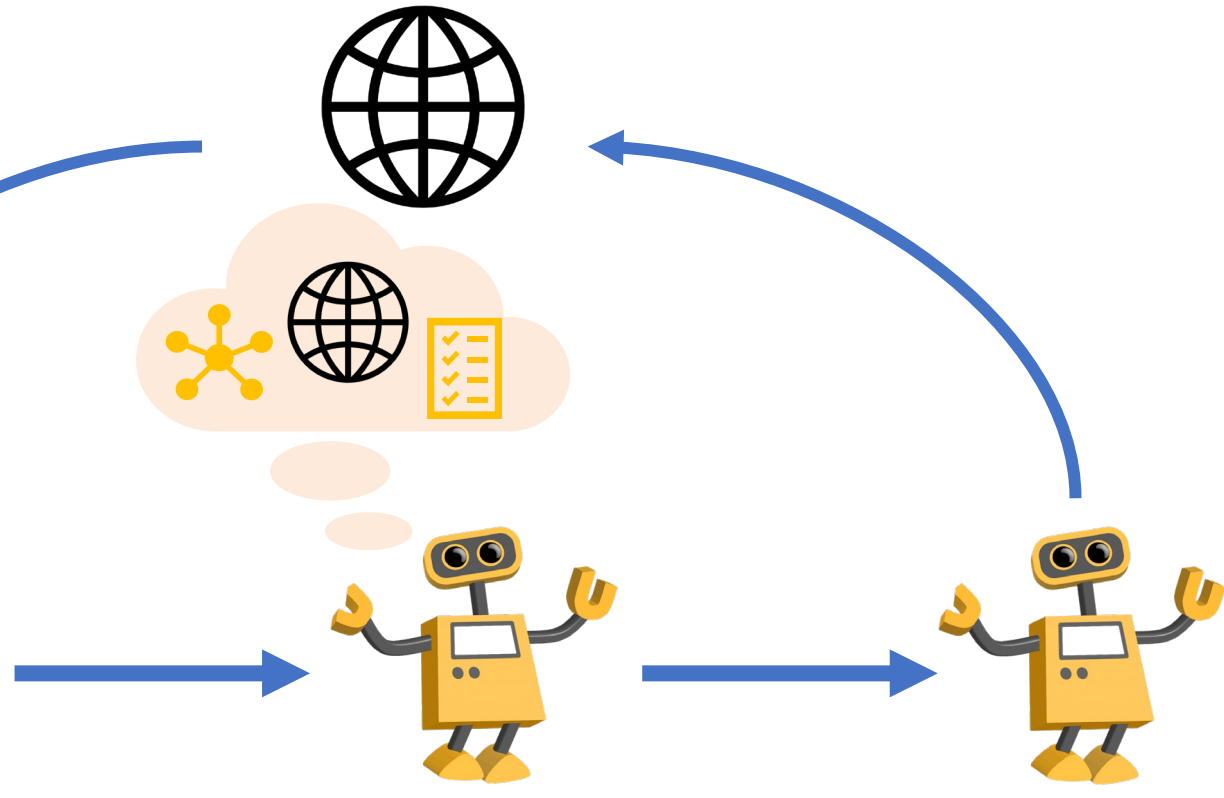
# Motivation: From Reactive to Deliberative

Agia et al., 2022  
Wang et al., 2022  
Srivastava et al., 2021  
She et al., 2014  
...

*Deliberative*



Situation



State Tracking  
& Planning

Action

# Contributions

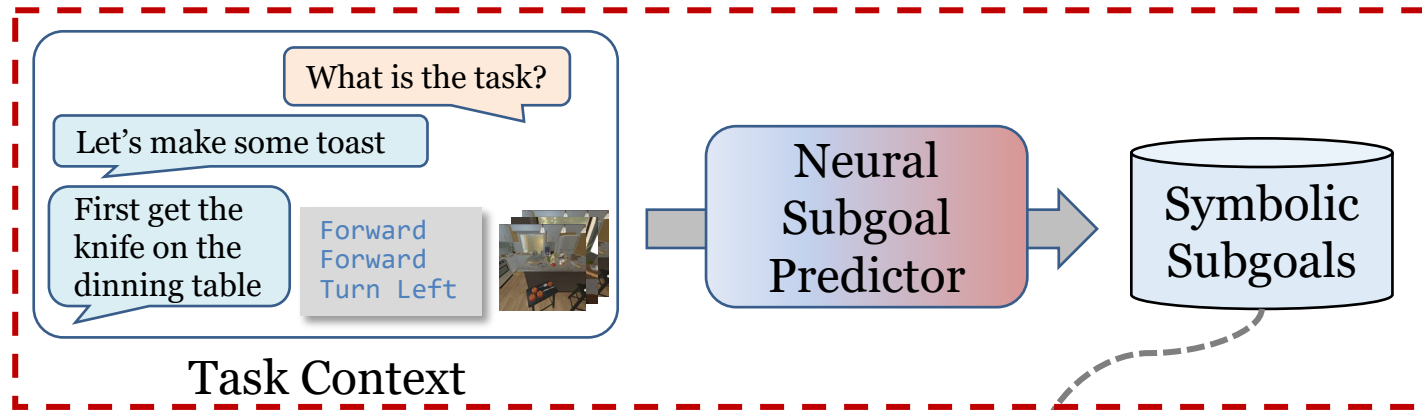
---

Deliberative Agent for following Natural Language Instructions (**DANLI**):

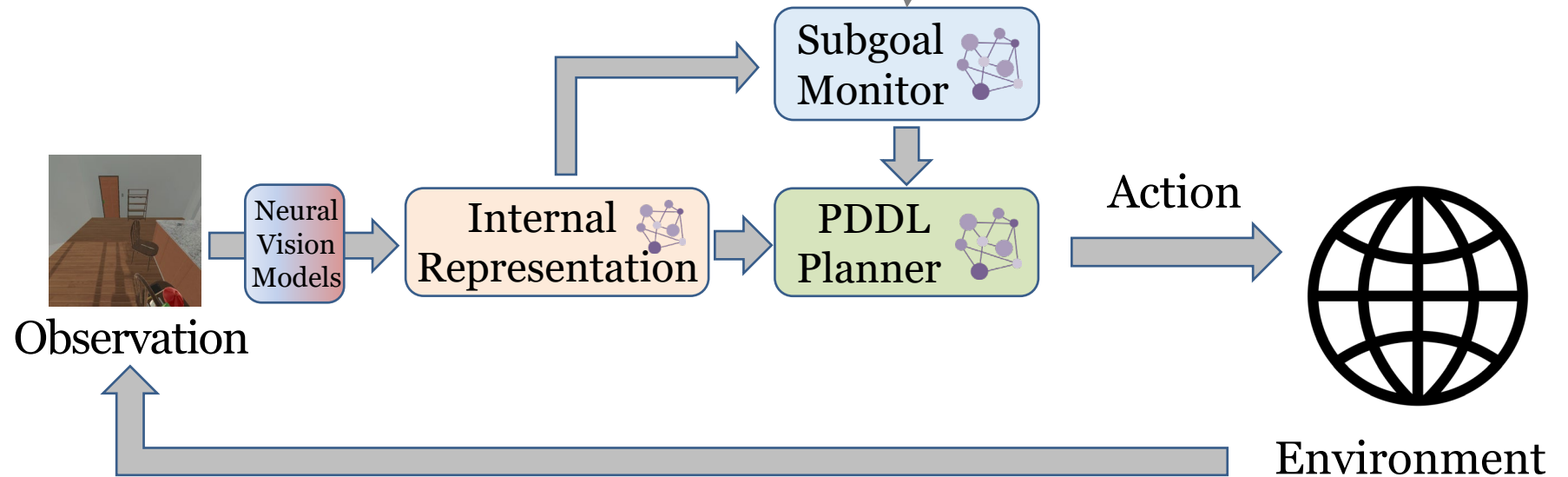
- 1. Neural subgoal predictor** uses a language model to predict symbolic subgoals from situated dialog history
- 2. Spatial-Symbolic world representation** tracks and grounds objects in a 3D voxel map to facilitate both path and task planning
- 3. Online symbolic planner** generates efficient, interpretable plans while allowing online exception handling

# Neuro-Symbolic System Architecture

## Subgoal Prediction



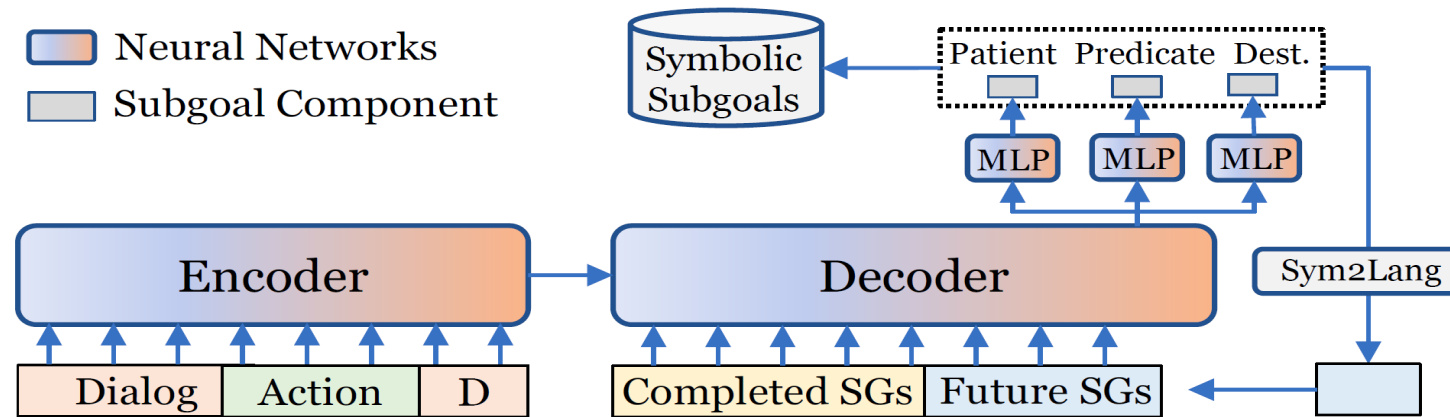
## Rollout





# Subgoal Learning From Dialog

- Neural encoder-decoder network for subgoal prediction from interaction history
- Joint progress estimation and future subgoal prediction



## Example Dialog & Action History Input

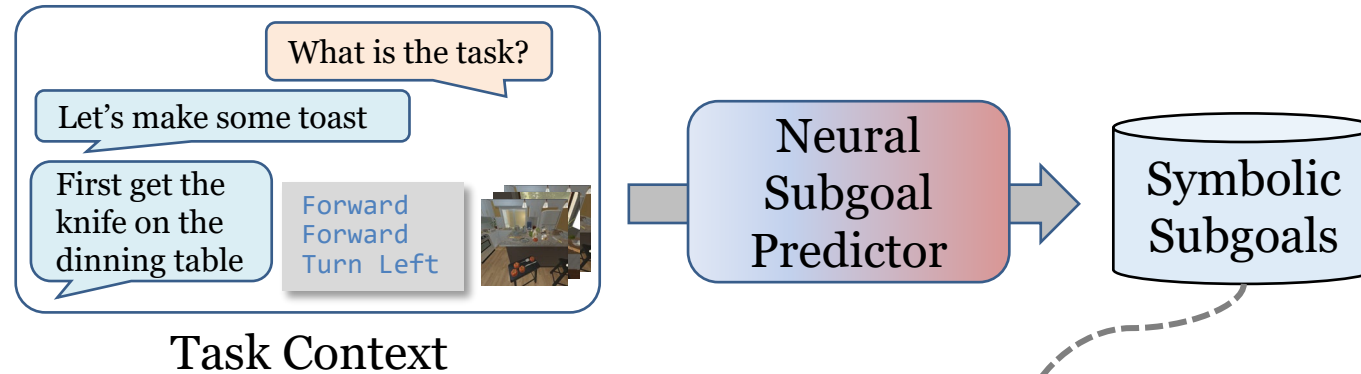
Follower: Hi. What can I do for you? Commander: Find a cup.  
Follower go for cup, pick up cup. Commander : Put it on the table.

## Example Completed & Future Subgoal Output

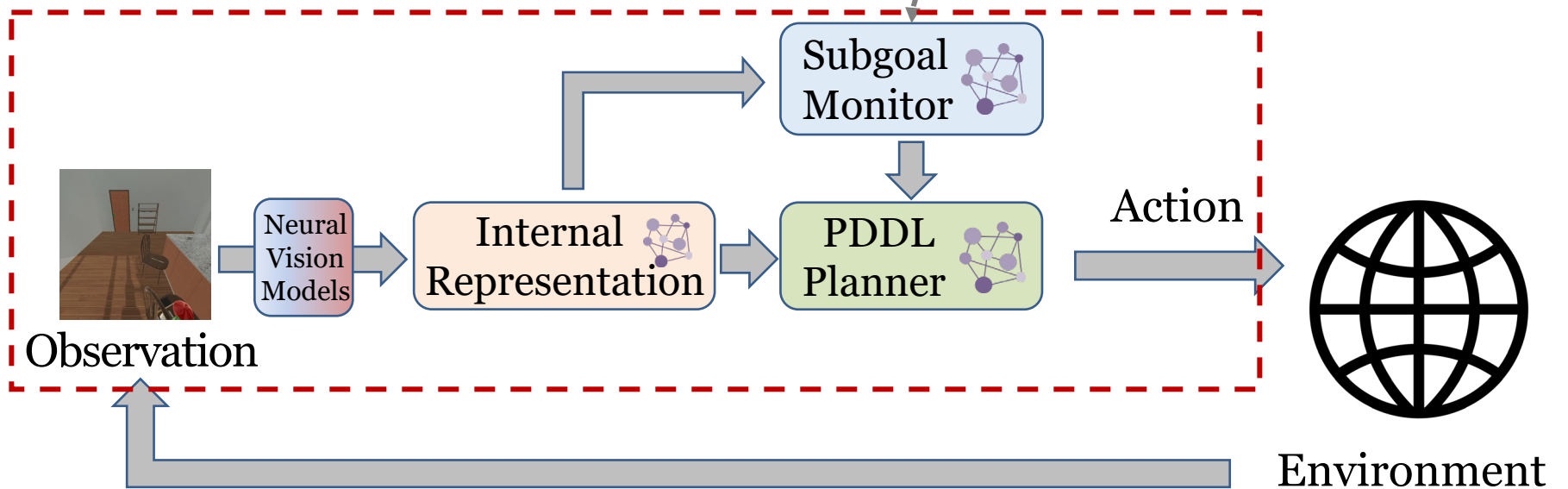
**Completed SG:** (Cup, isPickedUp)  
**Future SG:** (Cup, isPlacedTo, Table)

# Neuro-Symbolic System Architecture

## Subgoal Prediction

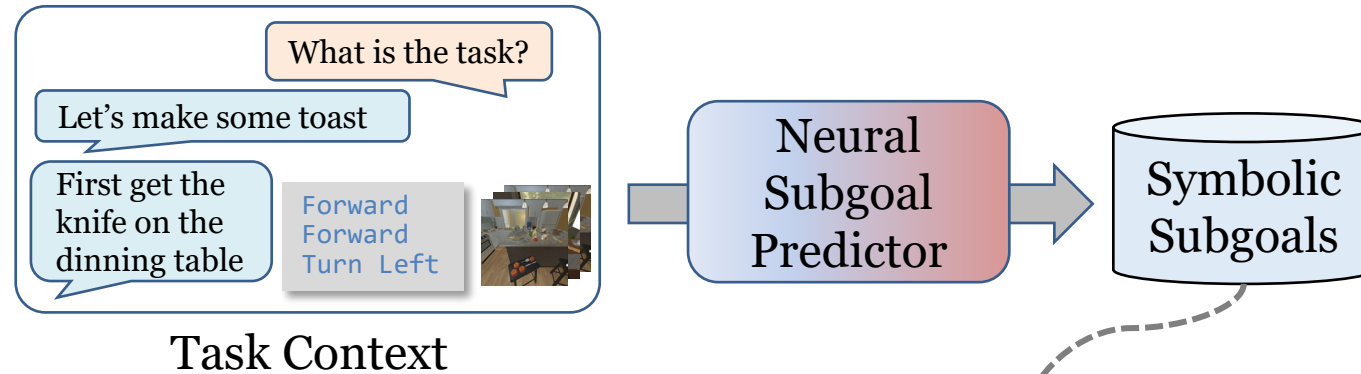


## Rollout

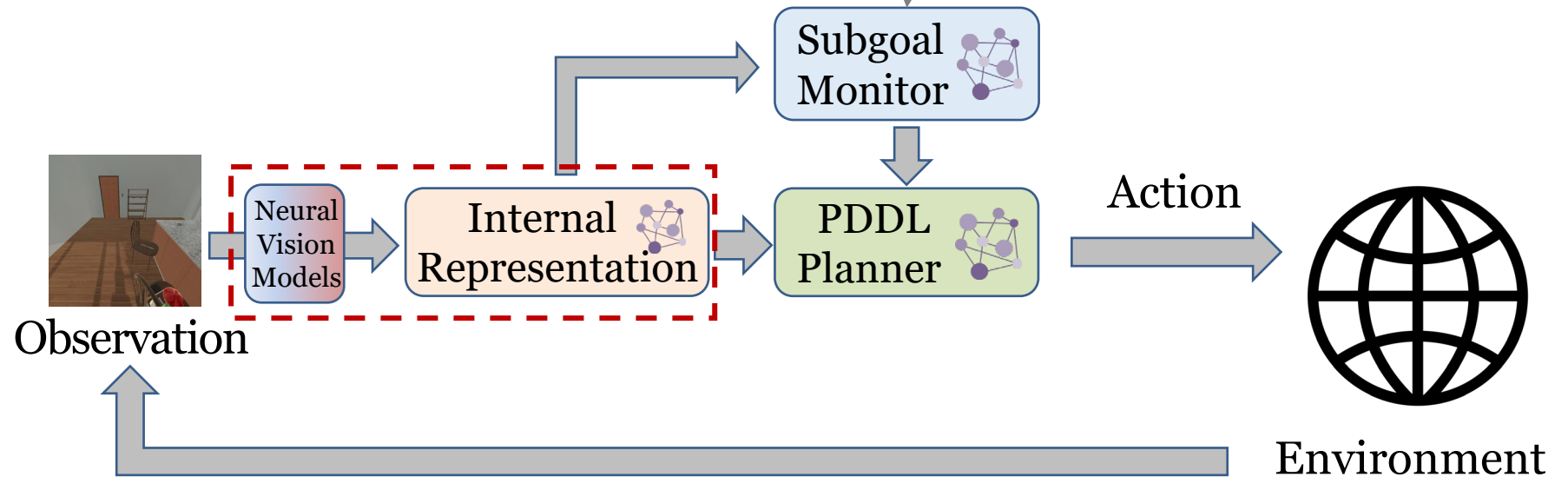


# Neuro-Symbolic System Architecture

## Subgoal Prediction

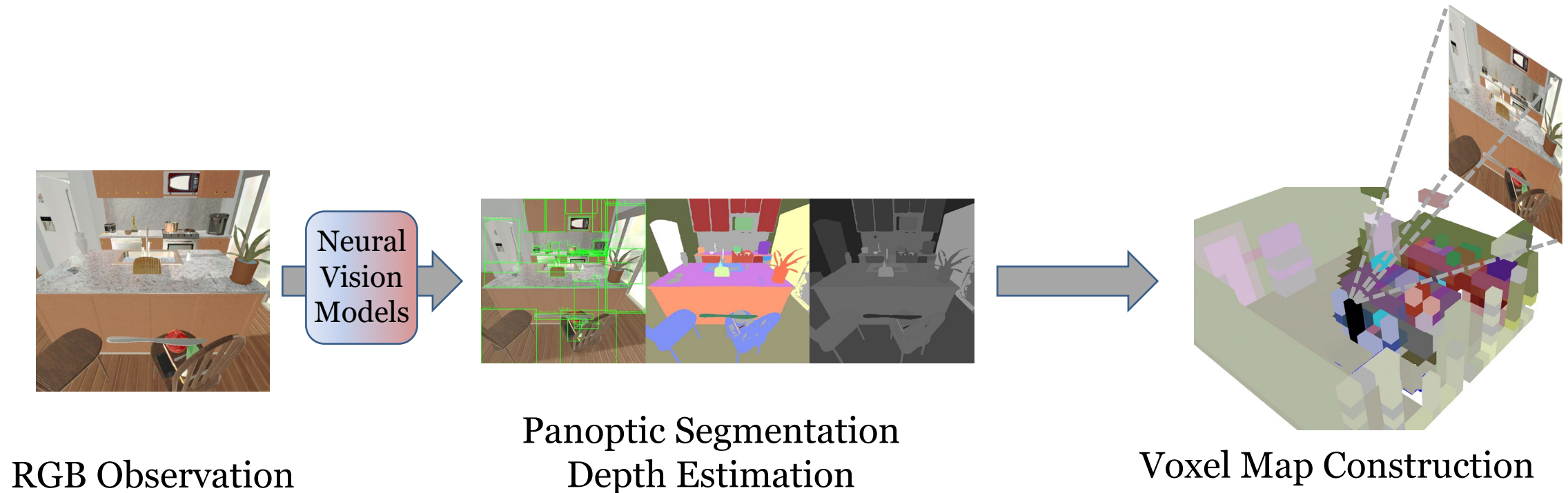


## Rollout



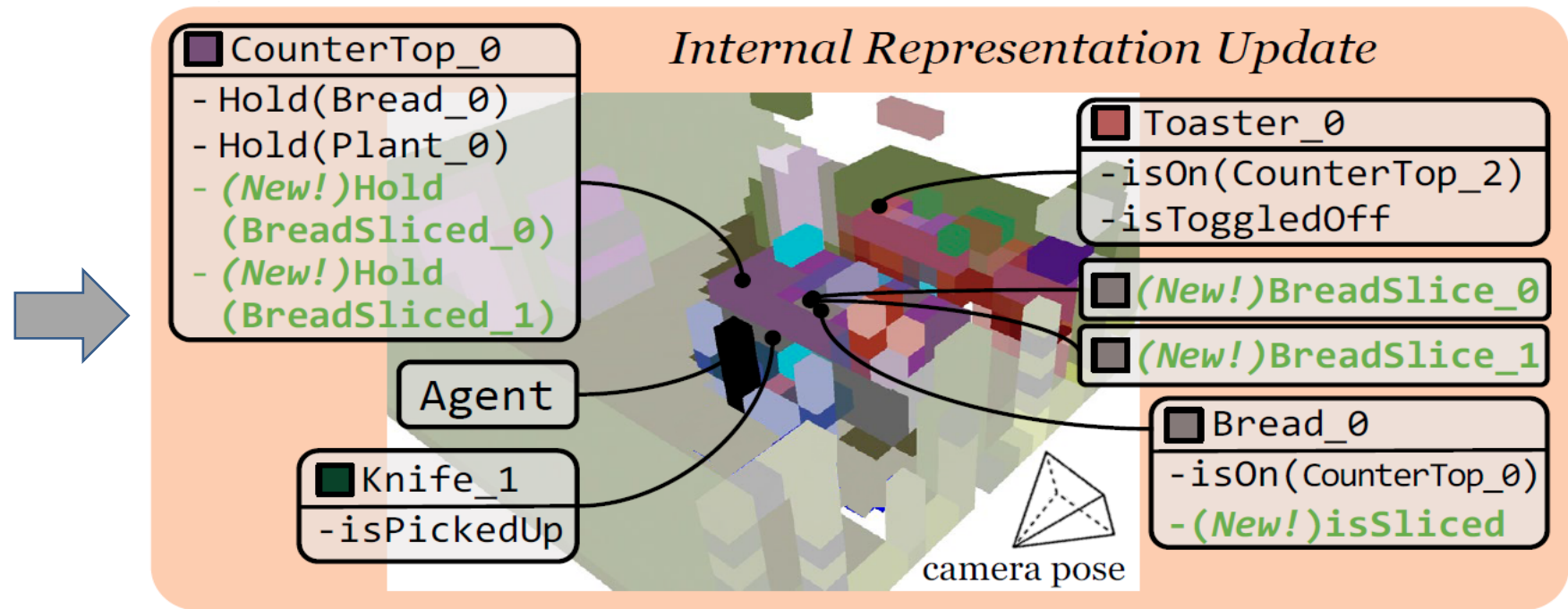
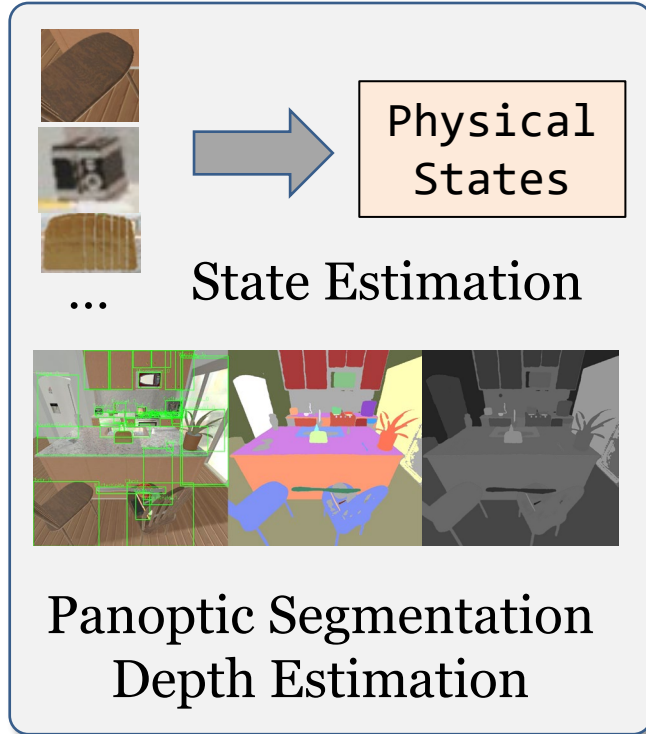
# Spatial-Symbolic World Representation

- Neural-powered 3D voxel map construction from ego-centric observations



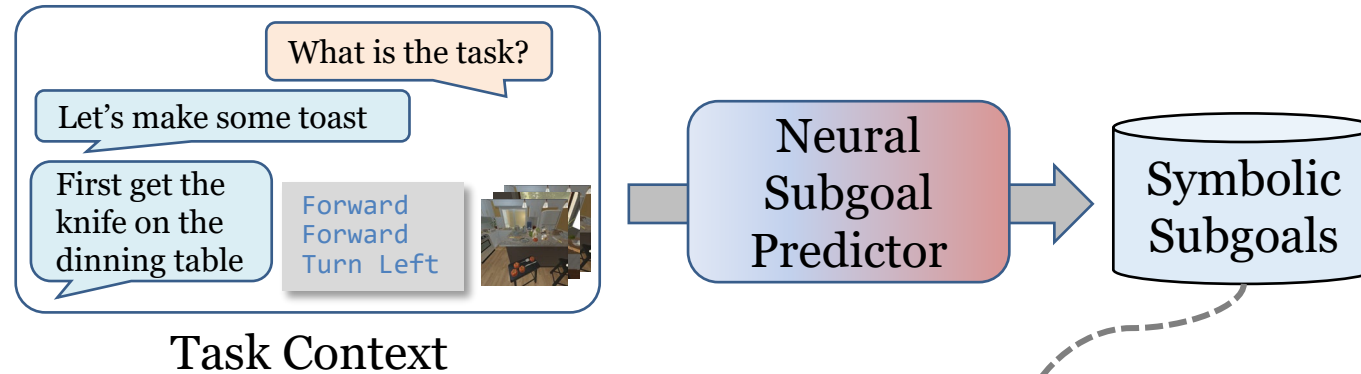
# Spatial-Symbolic World Representation

- Recognize physical states of each object instance
- Update the spatial-symbolic world representation at every step

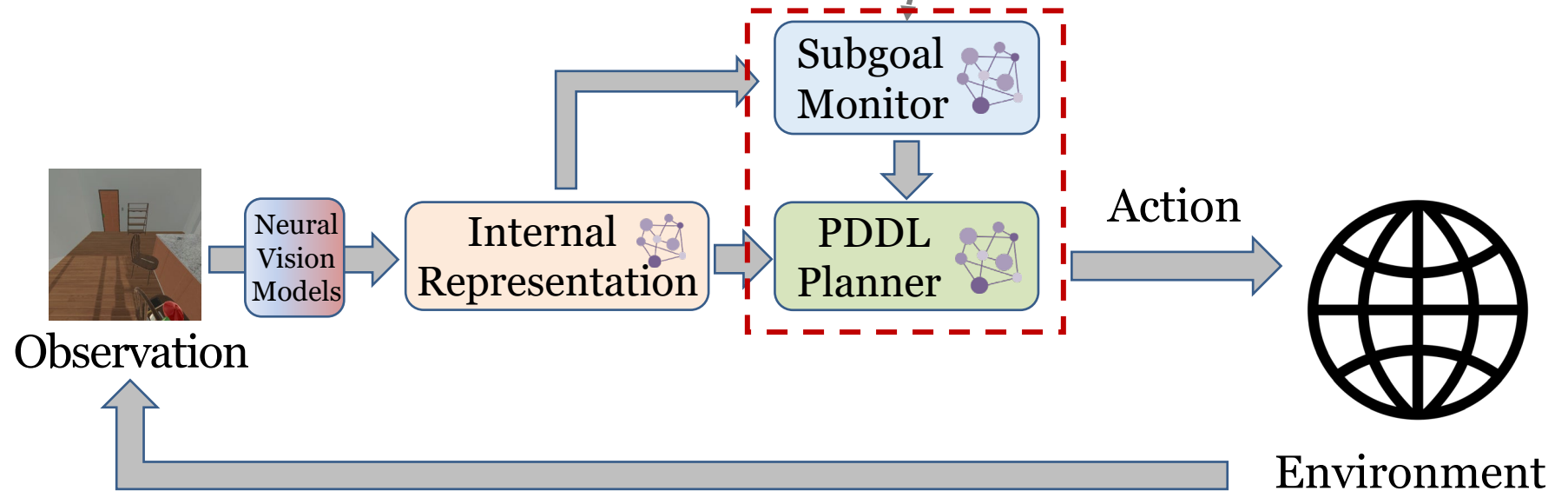


# Neuro-Symbolic System Architecture

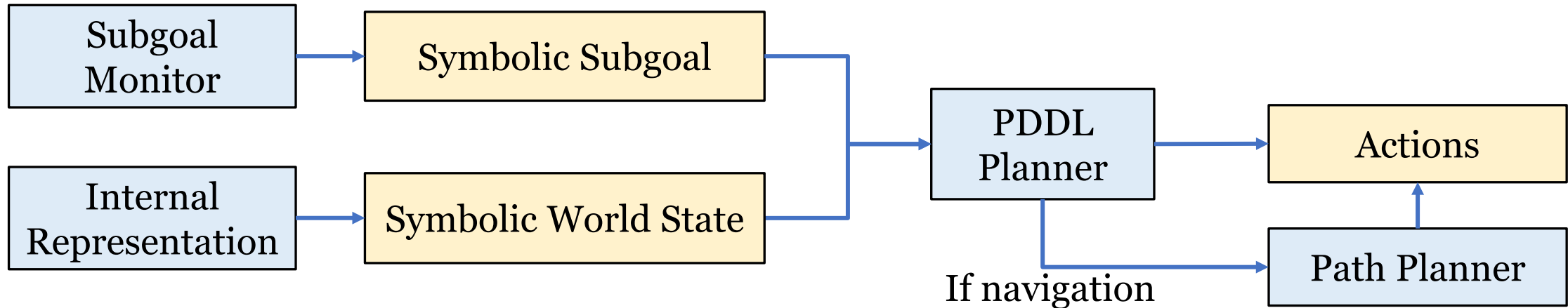
## Subgoal Prediction



## Rollout



# Symbolic Planning Pipeline



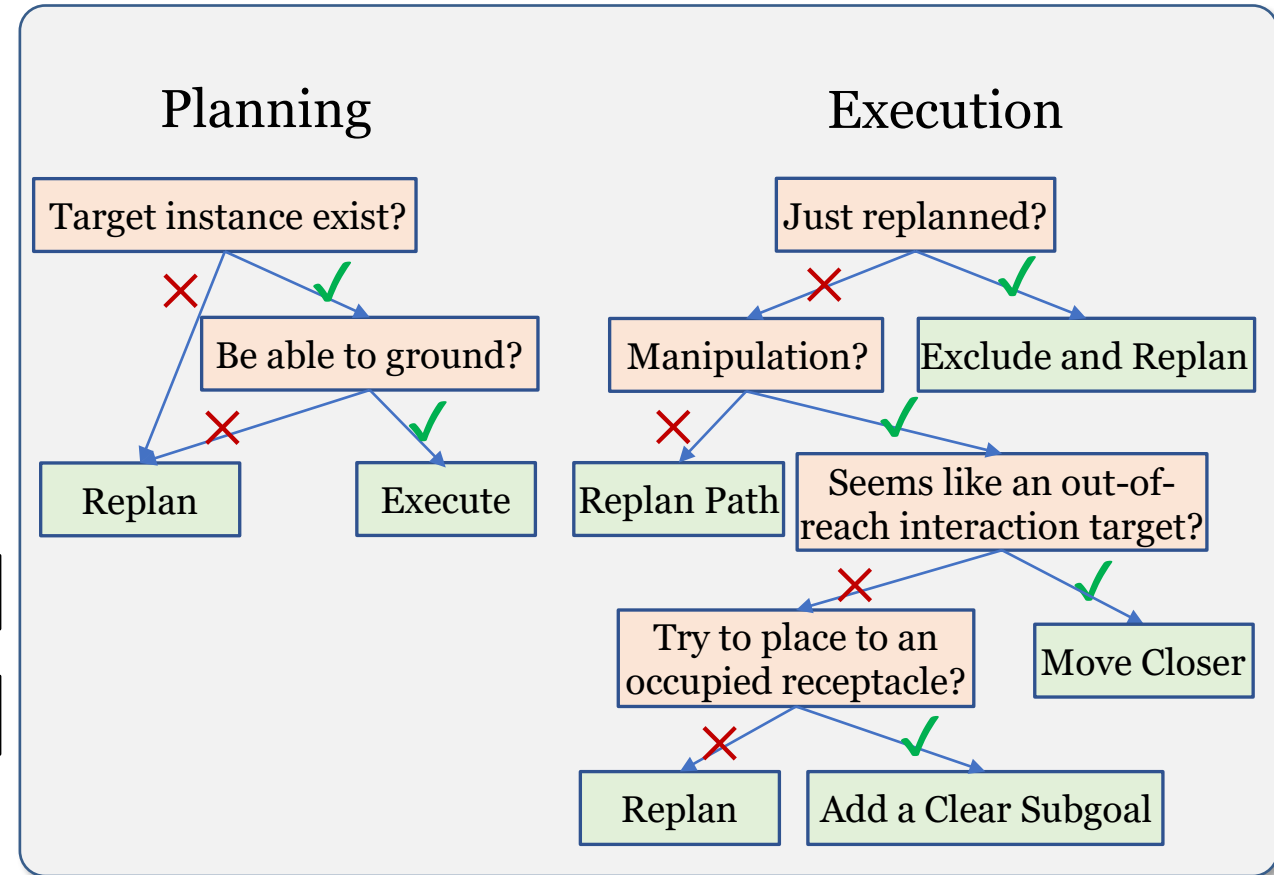
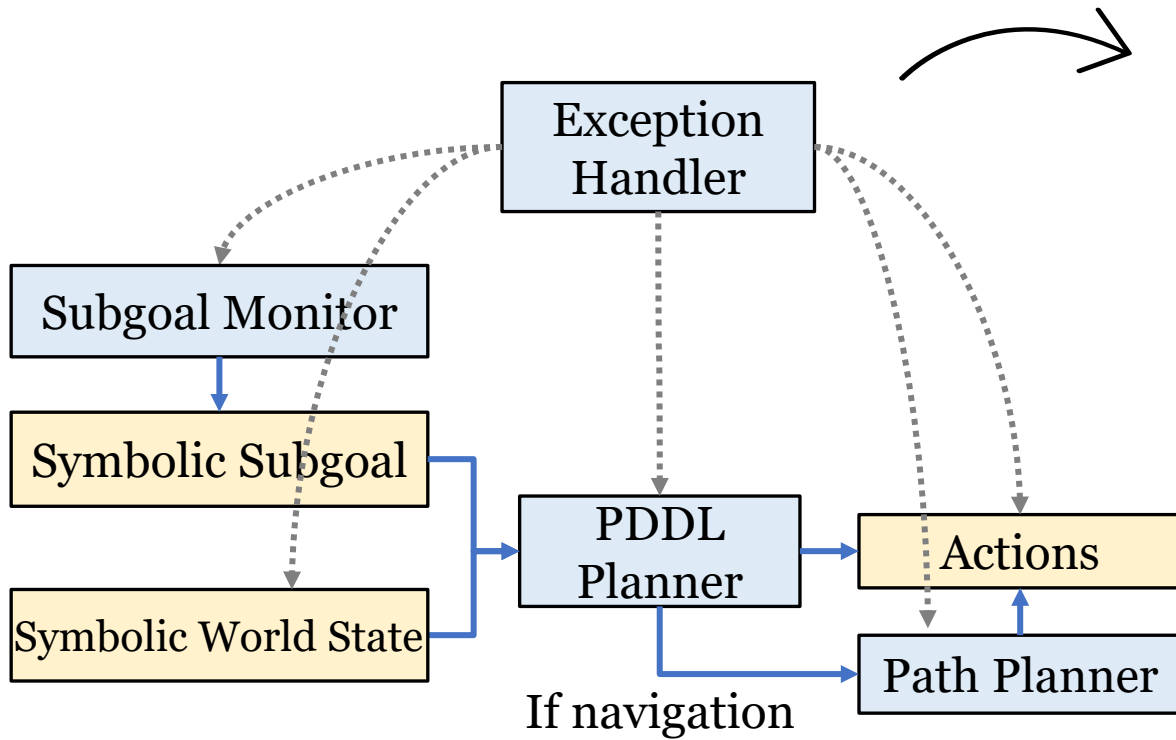
```
(:goal
  (and
    (exists (?o1 - BreadSliced ?o2 - Toaster)
      (and (parentReceptacles ?o1 ?o2))
    )
  )
)
```

Snapshot of an example subgoal  
“putting a piece of bread into toaster”

```
(:action Slice
  :parameters (?x - Sliceable ?y - Knives)
  :precondition (and (holding ?y) (isInteractable ?x) (not (isSliced ?x)) (not (isPickedUp ?x)))
  :effect (and
    (isSliced ?x)
    (not (isObserved ?x))
    (increase (total-cost) 15)
  )
)
```

Snapshot of the action “Slice” defined in the PDDL domain

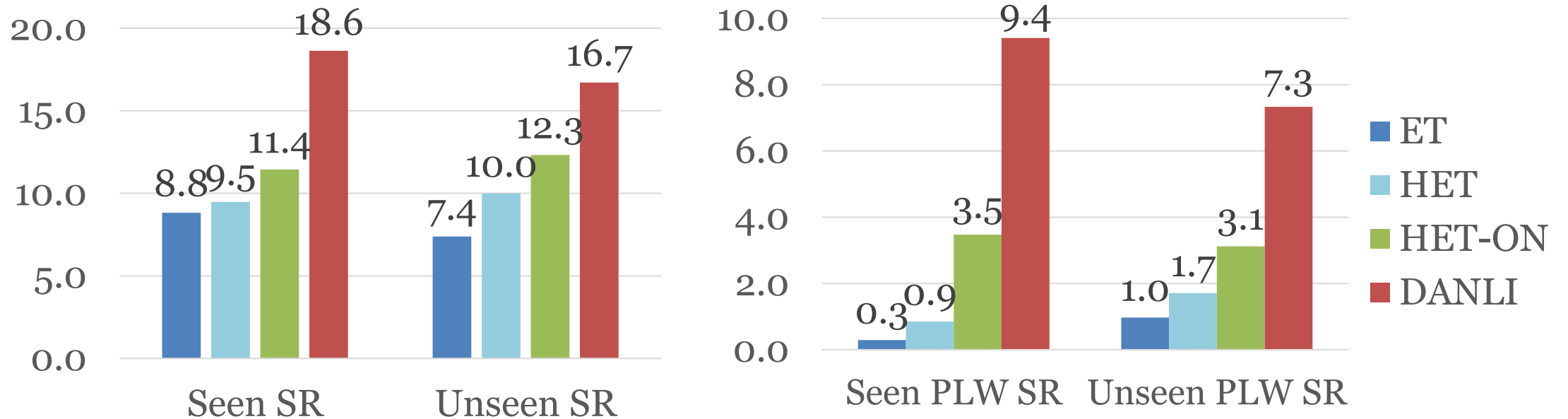
# Exception Handling In Online Planning





# Higher Success Rates

- DANLI outperforms reactive agents (ET, HET, HET-ON) by a large margin, especially when the task completion efficiency is considered

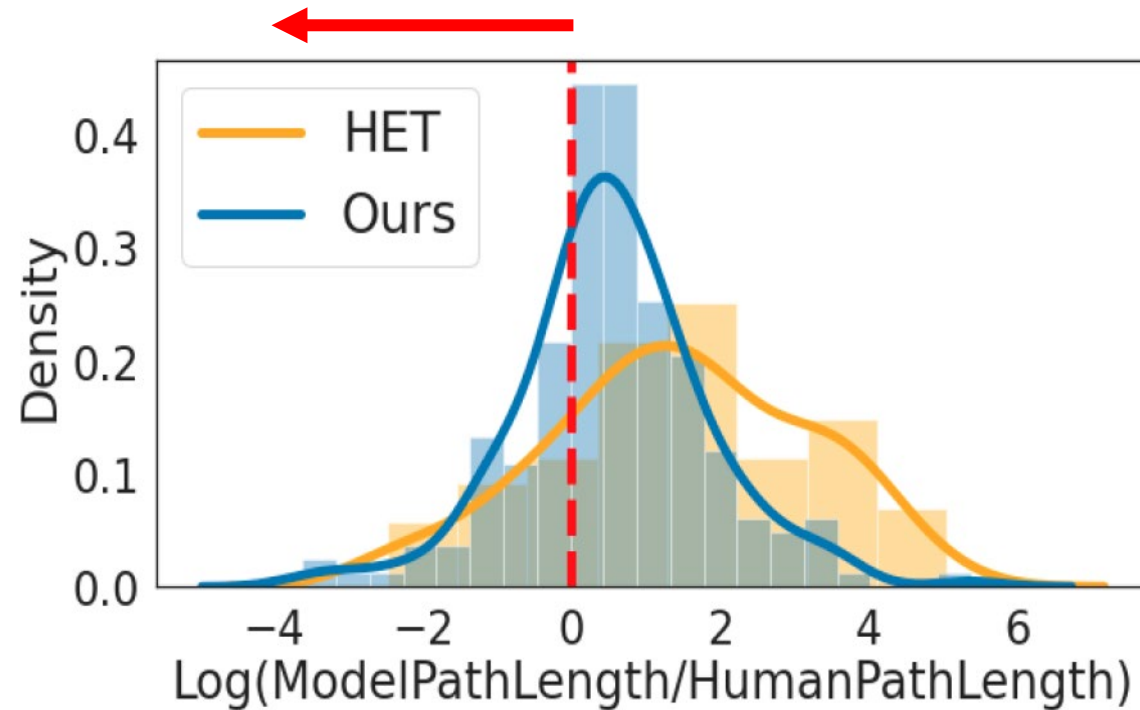


- Success Rate (SR): the proportion of successfully completed tasks
- Seen/Unseen: whether the evaluated scene is seen or unseen during training
- PLW SR: success rate weighted by the relative trajectory path length compared to humans

# Higher Efficiency

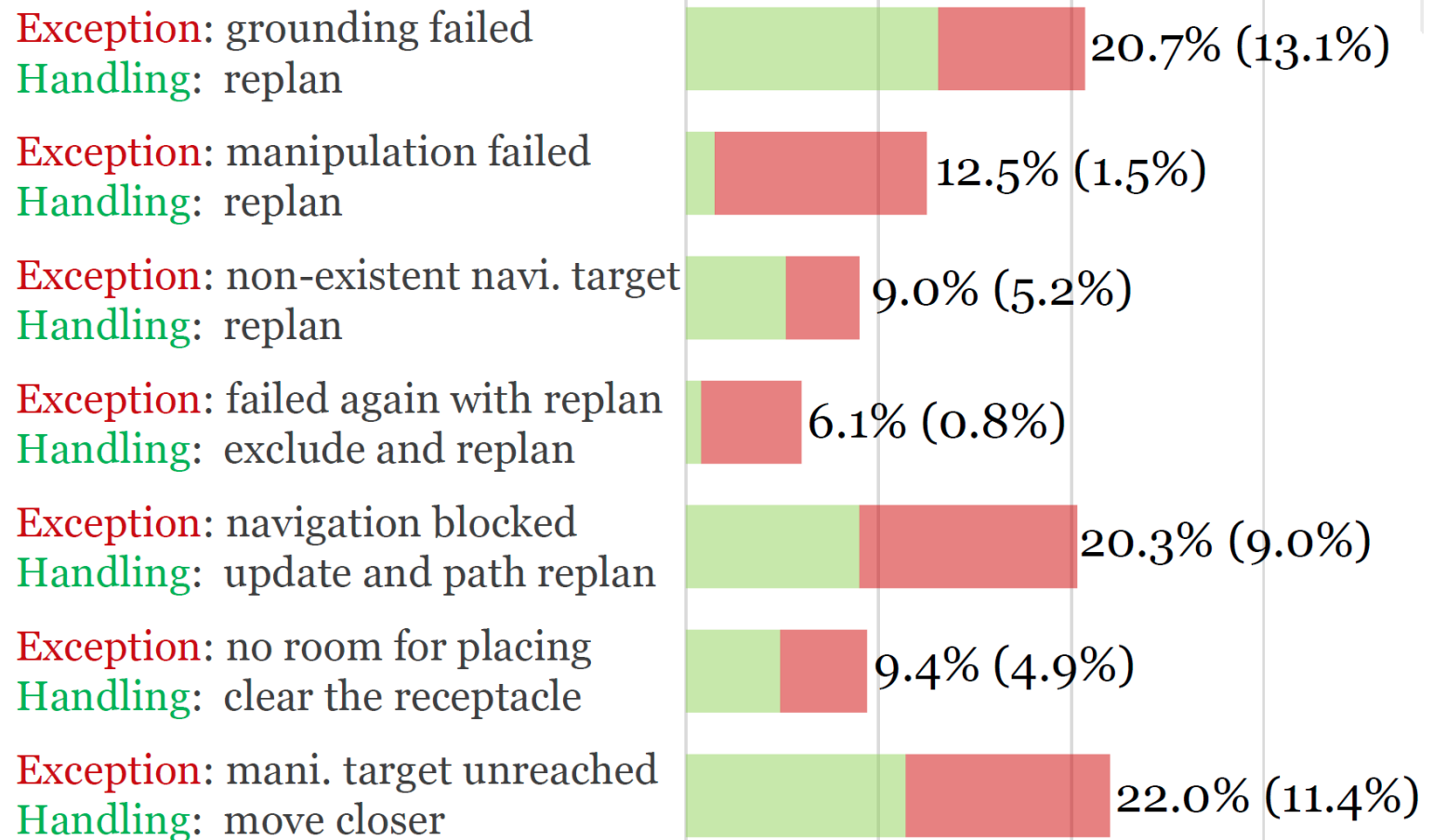
- DANLI completes tasks in fewer steps than reactive agents

Shorter path than human trajectories



# Interpretable Exception Handling

- Deliberative agent produces interpretable plans which allow fine-grain exception handling.



# Demo

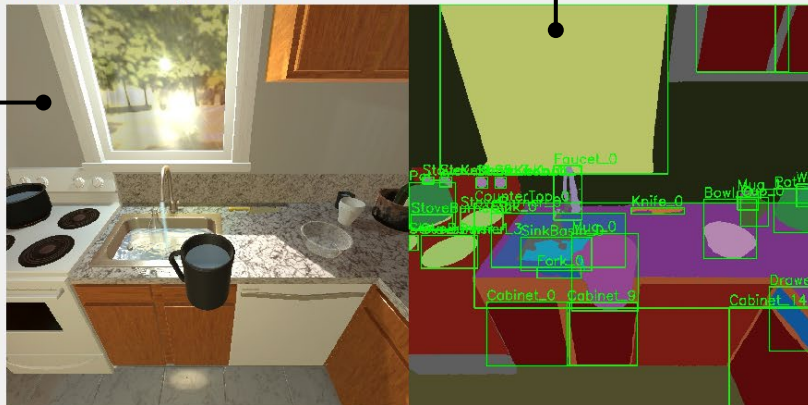
Panoptic Segmentation Output

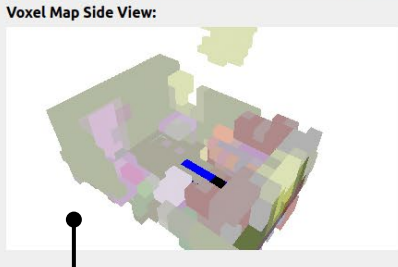
Step Number

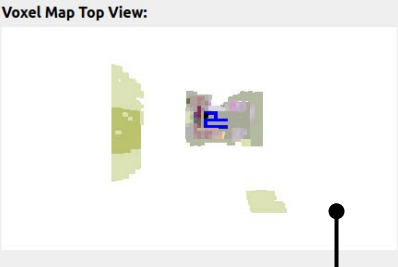
Stage: Replay/Rollout

Raw Observation

Raw Observation & Panoptic Segmentation:



Voxel Map Side View: 

Voxel Map Top View: 

Step: 1 Stage: rollout

Dialog:

BOT: what's first?  
USR: hi  
BOT: hey  
USR: prepare a coffee in clean mug  
BOT: where is the mug?  
USR: mug is right to the sink  
USR: good job

Events:

Get the initial plan for subgoal: Mug\_isEmptied

Subgoals:

- SG0: (Mug, isClean)
- SG1: (Mug, isEmptied)
- SG2: (Mug, simbotisFilledWithCoffee)

Plan for the Subgoal:

- A0: Pour(Bowl\_0)
- A1: Stop

Current Subgoal: SG1: (Mug, isEmptied)

Next Action: Pour(Bowl\_0)

Previous Next

Dialog History

Key Reasoning Events

Action Plan for the Current SG

Subgoals predicted by DANLI

Voxel Map constructed by DANLI  
Left: Side View / Right: Top-Down View

Current Subgoal

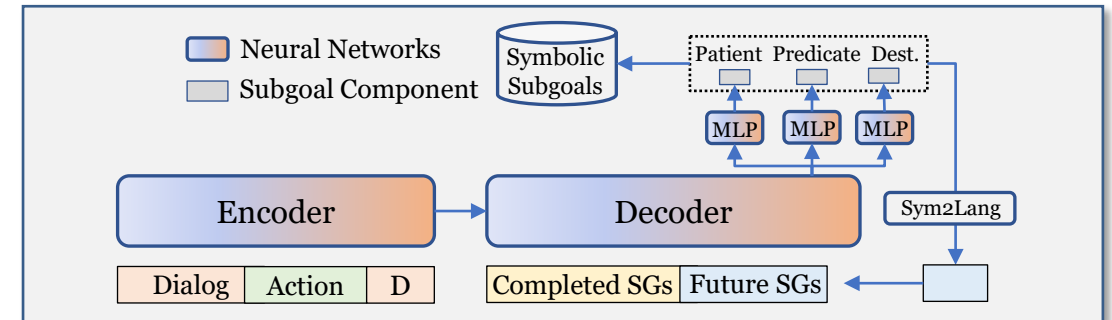
Next Action to Execute



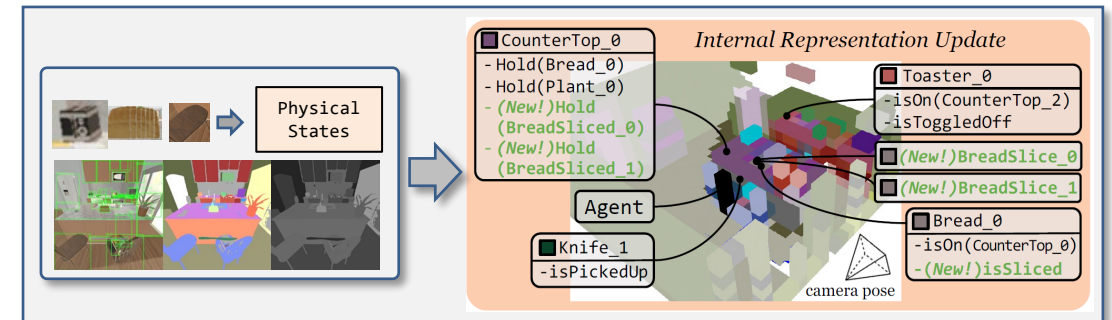
# Summary - DANLI



## Neural Subgoal Predictor



## Spatial-Symbolic World Representation



## Symbolic Planner with Online Exception Handling

