# Asking the Right Questions: Learning Interpretable Action Models Through Query Answering

**Pulkit Verma, Shashank Rao Marpally, Siddharth Srivastava,** Arizona State University
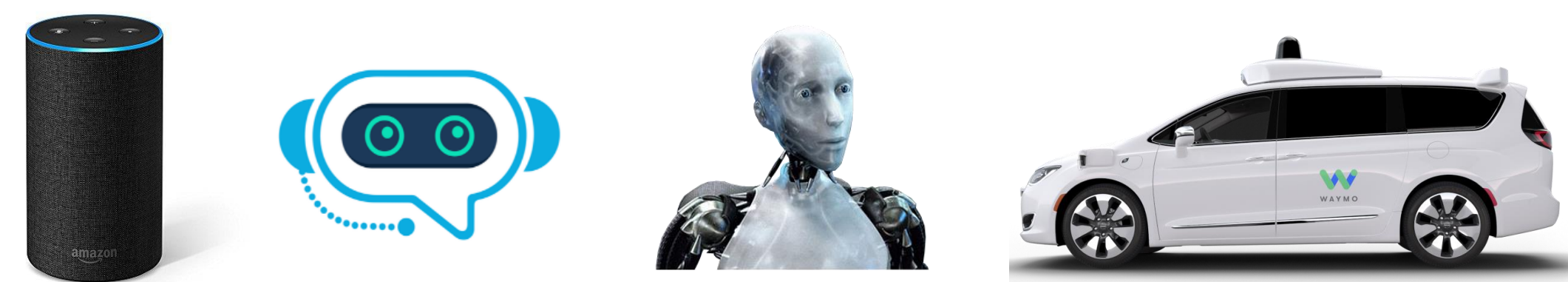
*How would a non-expert assess the limits and capabilities of an AI system?*

## Introduction

Objective: Learn an interpretable model of a black-box agent by interrogating it.

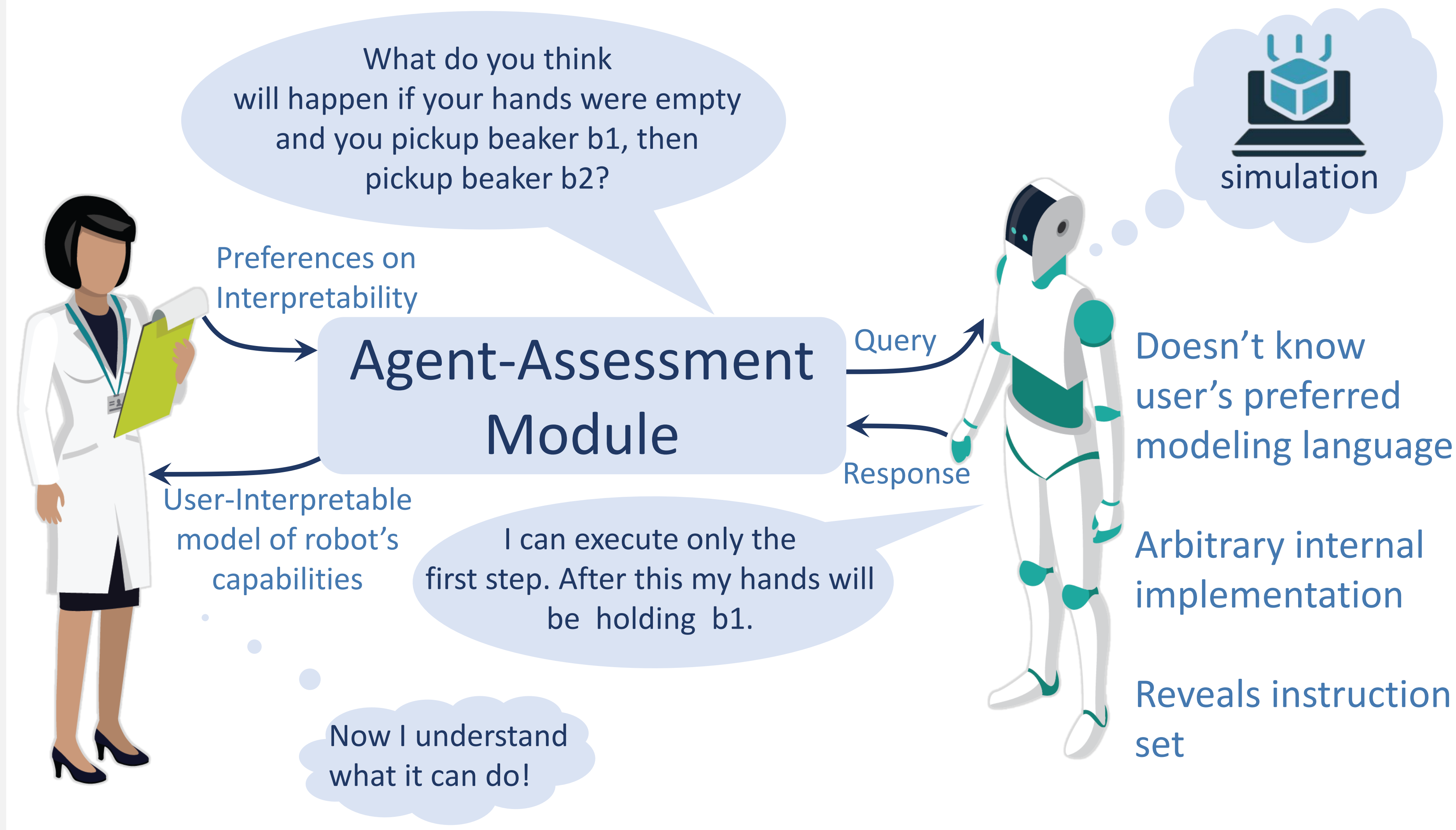Key technical challenge:
- Which sequence of queries to ask?

## Algorithm (AIA)

① Start with the most abstracted node in lattice.
② Pick abstraction candidates in some order.
③ For each candidate, generate three models and for each pair of models:
④ • Generate a distinguishing query $Q$ and pose it to the agent.
⑤ • Get the response R from the agent.
⑥ • Prune out the incorrect variants of candidate models.
⑦ • Repeat steps 3-6 till the model is fully estimated.
⑧ Return the final set of model(s).

## Salient Features

- Efficiently learns internal model of an autonomous agent in a STRIPS-like form.
- Needs no prior knowledge of the agent model.
- Only requires an agent to have rudimentary query answering capabilities.
- Queries can be answered using a simulator.

## Example of Agent Interrogation



What do you think will happen if your hands were empty and you pickup beaker b1, then pickup beaker b2?

Preferences on Interpretability

Agent-Assessment Module

User-Interpretable model of robot's capabilities

Query

Response

I can execute only the first step. After this my hands will be holding b1.

Now I understand what it can do!

simulation

Doesn't know user's preferred modeling language

Arbitrary internal implementation

Reveals instruction set

## Abstraction in Space of Models

```
(:action load_truck
 :parameters (?package ?truck ?location)
 :precondition (and (at ?truck ?location)
                    (+/-/∅) (at ?package ?location))
 :effect (and (not (at ?package ?location))
              (in ?package ?truck)))
```

**Abstracted model**

abstraction ↑

This predicate can appear in three forms:
- positive
- negative
- absent

```
(:action load_truck
 :parameters (?package ?truck ?location)
 :precondition (and (at ?truck ?location)
                    (at ?package ?location))
 :effect (and (not (at ?package ?location))
              (in ?package ?truck)))
```
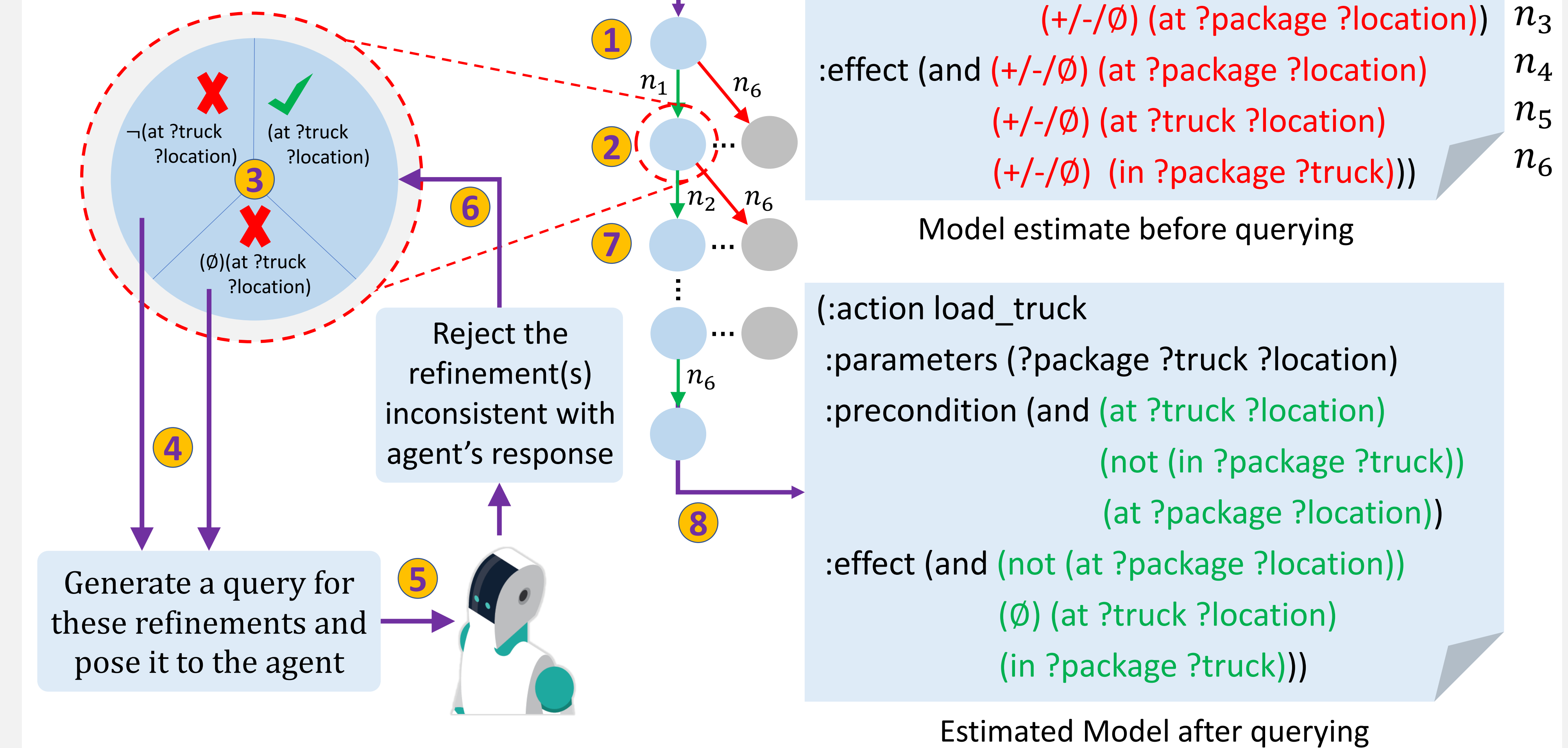
**Concrete model**
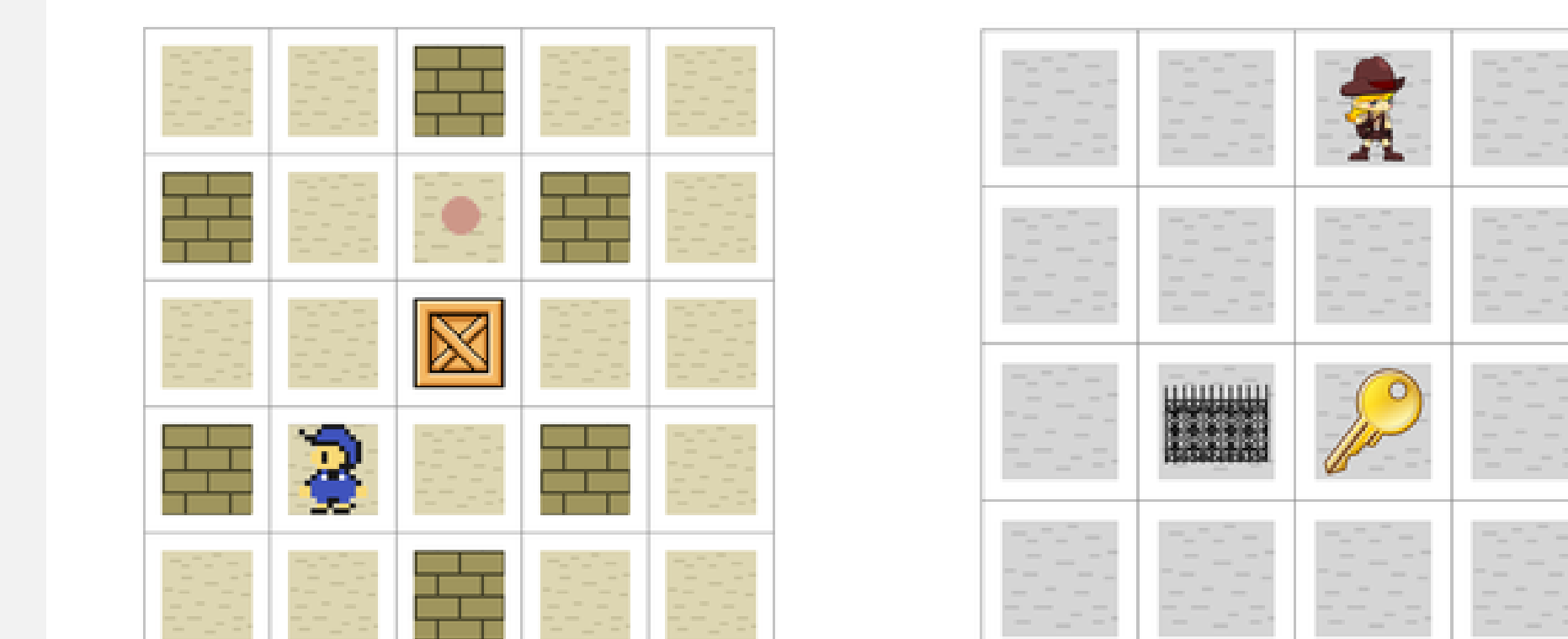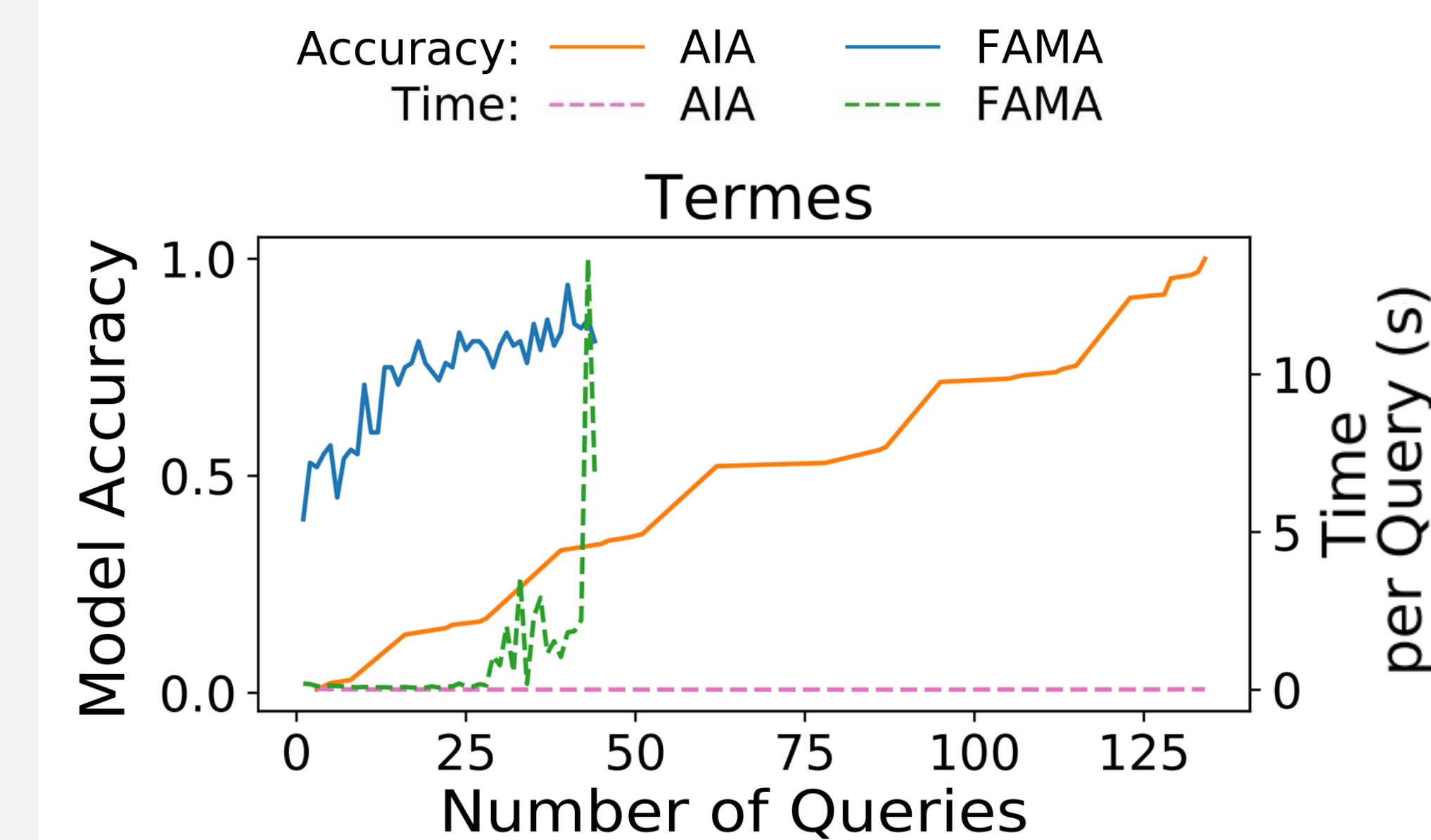
## Key Algorithmic Principle

**Key feature of the algorithm**
Each time we prune an abstracted model, we prune a very large number of models at the most concrete node.

Reject the refinement(s) inconsistent with agent's response

Generate a query for these refinements and pose it to the agent

```
(:action load_truck
 :parameters (?package ?truck ?location)
 :precondition (and (+/-/∅) (at ?truck ?location)      n_1
                    (+/-/∅) (in ?package ?truck)        n_2
                    (+/-/∅) (at ?package ?location))    n_3
 :effect (and (+/-/∅) (at ?package ?location)           n_4
              (+/-/∅) (at ?truck ?location)             n_5
              (+/-/∅) (in ?package ?truck)))            n_6
```
Model estimate before querying

```
(:action load_truck
 :parameters (?package ?truck ?location)
 :precondition (and (at ?truck ?location)
                    (not (in ?package ?truck))
                    (at ?package ?location))
 :effect (and (not (at ?package ?location))
              (∅) (at ?truck ?location)
              (in ?package ?truck)))
```
Estimated Model after querying

## Results



Accuracy: AIA   FAMA
Time: AIA   FAMA

Termes

- AIA efficiently derives interpretable agent models for a range of agents.
- AIA is much faster than state of the art methods for deriving models based on passive observations.
- AIA offers better convergence guarantees.

Refer to the paper for detailed results

bit.ly/3p4cVRu