

# User-Aligned Autonomous Capability Assessment of Black-Box AI Systems

#### Pulkit Verma and Siddharth Srivastava



# Focus: Taskable AI Systems

- Sequential Decision-Making Systems.
- User gives AI systems a task (an objective).
- Al system needs to figure out how to complete that task.



### Personalized Assessment of Adaptive AI Systems

- Users would like to give AI systems multiple tasks.
  - How would users know what the AI systems can do?
- Al systems should support third-party assessment.
- The assessment should work with black-box AI systems.



# Capability vs Functionality





#### Easier to reason in terms of capabilities than in terms of functionalities









Arbitrary internal implementation

## Black-Box AI System Interface



Personalized AI-Assessment Module • Should work for a variety of taskable AI systems

• Should be easy to support

Query	The plan $\pi = \langle c_1,, c_n \rangle$ ?	State $s_F$ ?
Response	I can execute first $\ell$ steps of the plan, ending up in state $s_F$ .	Yes / No.
	Plan Outcome Queries	State Reachability Query

at(p0,cell\_6\_3)
clear(cell\_0\_0)
door\_at(cell\_9\_2)
next\_to(m0)
alive(m0)
key\_at(9\_4)

[Input] Concepts that the user understands

Query Response Personalized Al-Assessment Module



Arbitrary internal implementation

Doesn't know user's vocabulary

Black-Box AI



#### Interpretable Description: PDDL/PPDDL

```
(:action open-door
  :parameters (?l1)
  :precondition (and
     (has_key)
     (player_at ?l1)
     (door_adjacent ?l1))
  :effect (probabilistic —
    0.95 (and (door_open))
     0.05 (and (not (has_key))
               (game-over))
```

Precondition: This condition must be true for this action to execute

Effect: This is a set of conditions, one of which becomes true when this action is executed

Probabilities: Each set of effect has an associated probability with which that effect set is executed

### Interpretable: Easily Convertible to Natural Language

```
(:action open-door
  :parameters (?l1)
  :precondition (and
     (has_key)
     (player_at ?l1)
     (door_adjacent ?l1))
  :effect (probabilistic
     0.95 (and (door_open))
     0.05 (and (not(has_key))
               (game-over))
```

The player can open the door when in location ?l1 if:

- It has the key
- The player is at location ?l1
- The door is adjacent to location ?l1 After executing that capability:
- With 95% probability, the door will open
- With 5% probability, the player will not have the key and the game will be over

## **Exponential** Search for Learning Correct Description

- Consider the following 4 predicates/concepts:
  - (has\_key)
  - (door\_open)
  - (door\_adjacent ?x)
  - (player\_at ?x)
- Consider just one capability: (open-door ?x)
- 9<sup>|C|×|P|</sup> = 9<sup>1×4</sup>=6561 possible models (Assuming deterministic models/ descriptions, i.e., no probabilities).

```
(:action open-door
  :parameters (?l1)
  :precondition (and
     (+/-/\emptyset) (has_key)
     (+/-/\emptyset) (door_open)
     (+/-/Ø)(door_adjacent ?l1)
     (+/-/\emptyset) (player_at ?l1))
  :effect (and
     (+/-/\emptyset) (has_key)
     (+/-/\emptyset) (door_open)
     (+/-/Ø) (door_adjacent ?l1)
     (+/-/\emptyset) (player_at ?l1))
```















#### Key feature of the algorithm

Whenever we prune an abstract model, we prune a large number of concrete models.

#### **Active Learning**

[Verma, Marpally, Srivastava; AAAI '21]

### **Evaluation with Known Capabilities**



#### Assumptions

- User's vocabulary matches simulator's vocabulary.
- Black-Box AI provides a list of capabilities.
- Deterministic setting.
- Randomly generate an agent and environment from International Planning Competition (IPC).
- Evaluate performance of the assessment module and compare it with FAMA<sup>†</sup>.

#### AAM learns Accurate Model with fewer Queries

- Random traces as input to FAMA.
- Increased #traces till it ran out of memory.

Accuracy: <u>AAM</u> — FAMA Time: <u>AAM</u> ---- FAMA





#### **Capability Discovery**

[Verma, Marpally, Srivastava; KR '22]



#### **Differential Assessment**



#### **Stochastic Setting**



#### **Capability Discovery**



#### **Differential Assessment**

[Nayyar\*, Verma\*, Srivastava; AAAI '22]



#### **Stochastic Setting**



#### **Capability Discovery**



#### **Differential Assessment**





#### **Stochastic Setting**

[Verma, Karia, Srivastava; NeurlPS '23]

User-Aligned Autonomous Capability Assessment of Black-Box AI Systems Pulkit Verma and Siddharth Srivastava

