

Foundation Models Meet Embodied Agents



Manling Li
Northwestern



Yunzhu Li
Columbia



Jiayuan Mao
MIT



Wenlong Huang
Stanford



Northwestern
University



COLUMBIA



Stanford
University

Part I: Motivation and Overview

Manling Li, Assistant Professor at Northwestern University

Tutorial: Foundation Models Meet Embodied Agents



Northwestern
University



COLUMBIA



Stanford
University

What is a generalist agent?

What is a generalist agent?



Having a robot that can do many tasks, across many environments.

BEHAVIOR-1K



simulating and benchmarking robot tasks that **matter** to humans

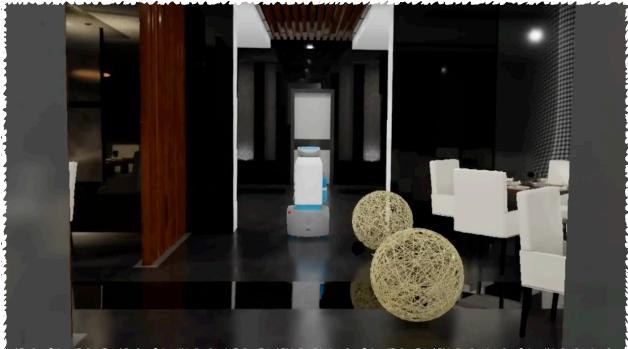
<https://behavior.stanford.edu/>

What is a generalist agent?



Having a robot that can do many tasks, across many environments.

Observation



State: 3D assets & states



<https://behavior.stanford.edu/>

Transition Model



What is a generalist agent?



Having a robot that can do many tasks, across many environments.

tasks that **matter**

What would you like a robot to help you with?



Cleaning the floor?

Images Generated by DALL-E 3

tasks that **matter**

What would you like a robot to help you with?



Folding Laundry?



Images Generated by DALL-E 3

tasks that **matter**

What would you like a robot to help you with?



Cooking Breakfast?



Images Generated by DALL-E 3

tasks that **matter**

What would you like a robot to help you with?



Opening gifts?



Images Generated by DALL-E 3

tasks that **matter**

What would you like a robot to help you with?



Cleaning after a wild party?

Images Generated by DALL-E 3

Robotic Manipulation



Human: pack the yoshi figure in the brown box



LEO:

$$P = [0.1, -0.2, 0]$$

$$R = [0, 0, 0, 1]$$

Task Planning



Human: Organize and tidy up the bedroom.

LEO:

1. *Clean the floor by sweeping to remove ...*
2. *Make the bed ...*
3. *Fold and organize ...*

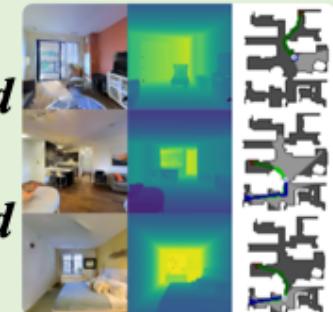
Embodied Navigation



Human: Find the pillow

LEO:

1. *Move forward*
2. *Turn right*
3. *Move forward*
4. ...



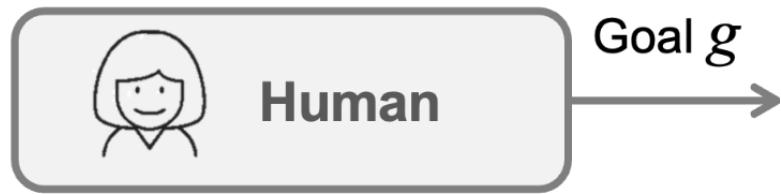
To achieve these tasks in these environments

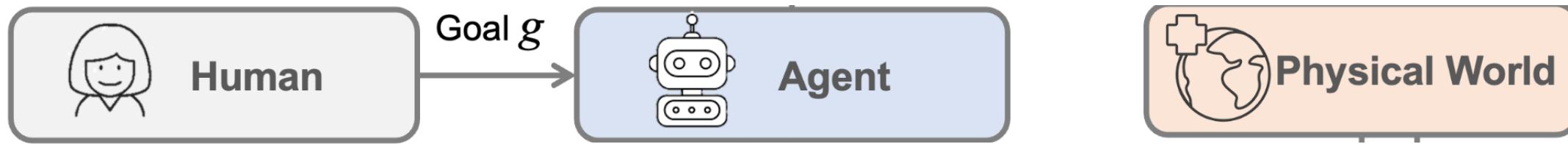
“embodied decision making”

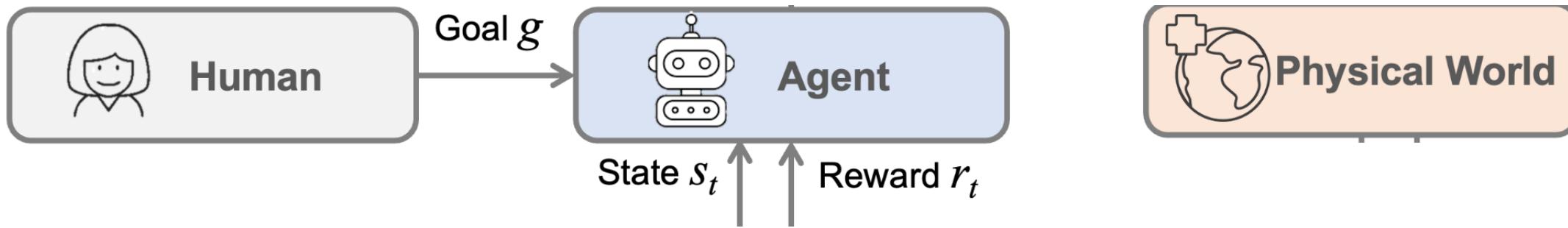
What is “embodied decision making” ?

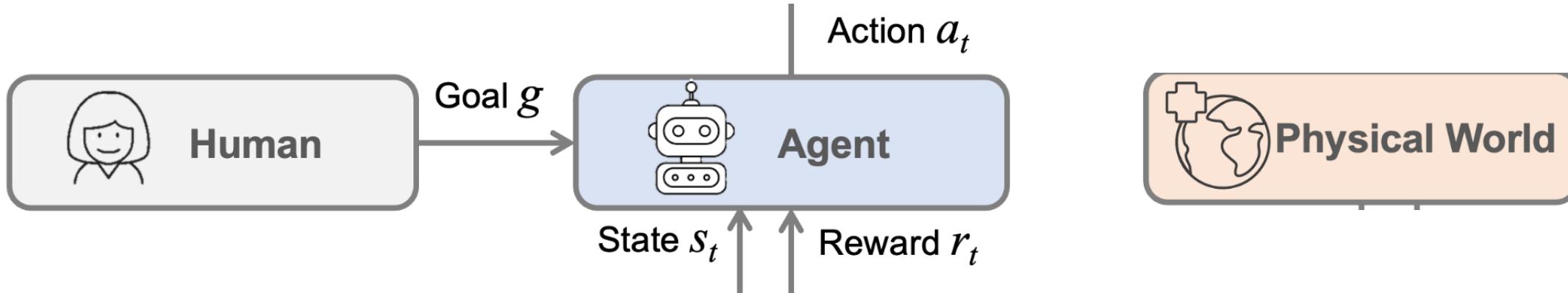
Can we leverage MDP as a guiding principle to categorize “foundation models”?

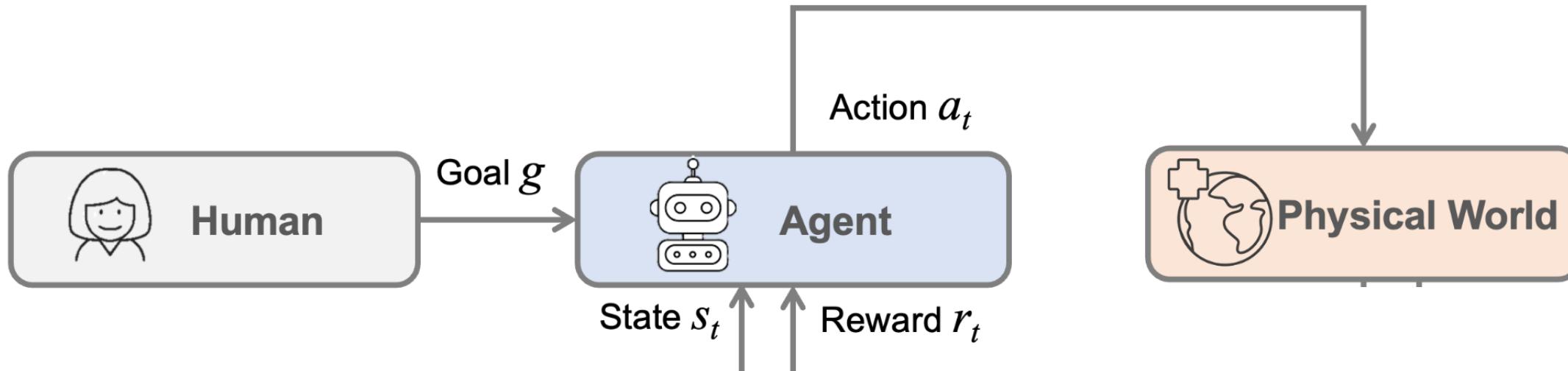
Let us go back to MDPs (Markov Decision Processes)

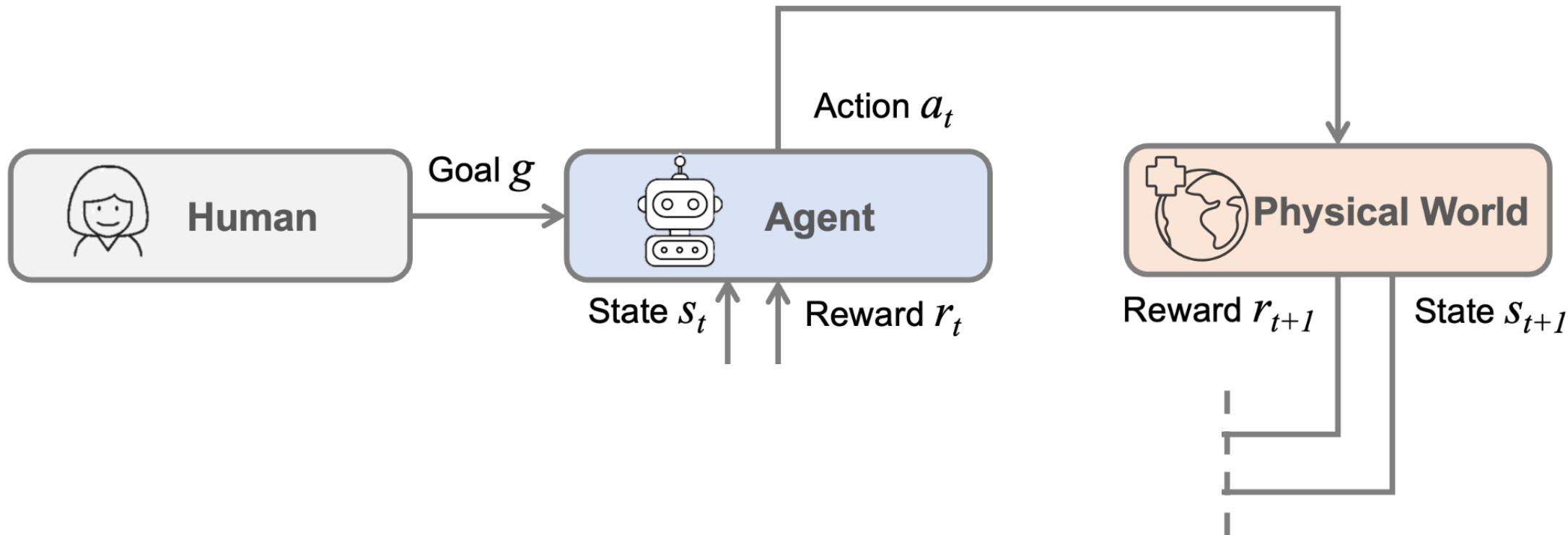


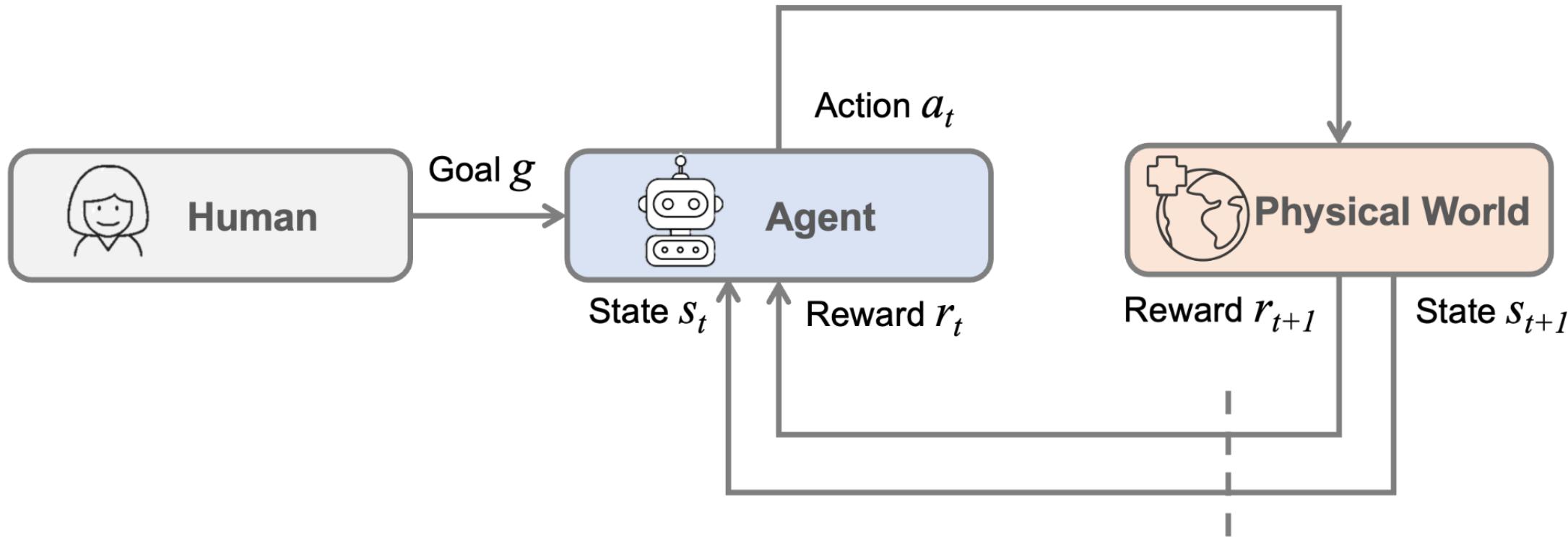


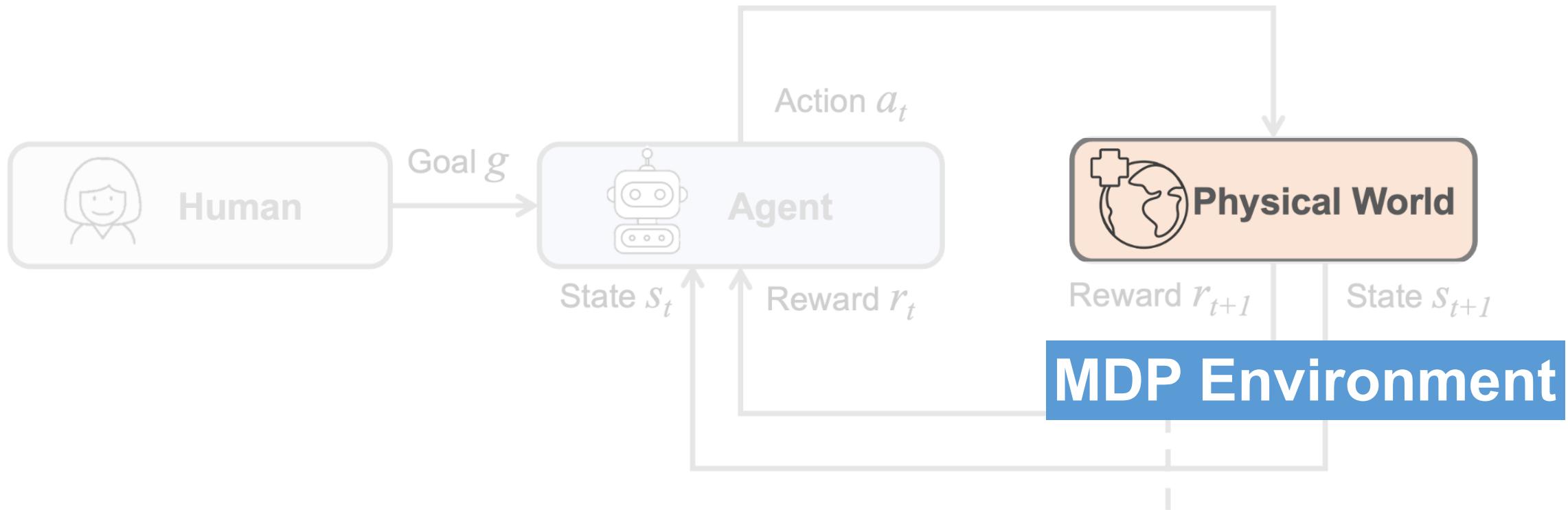


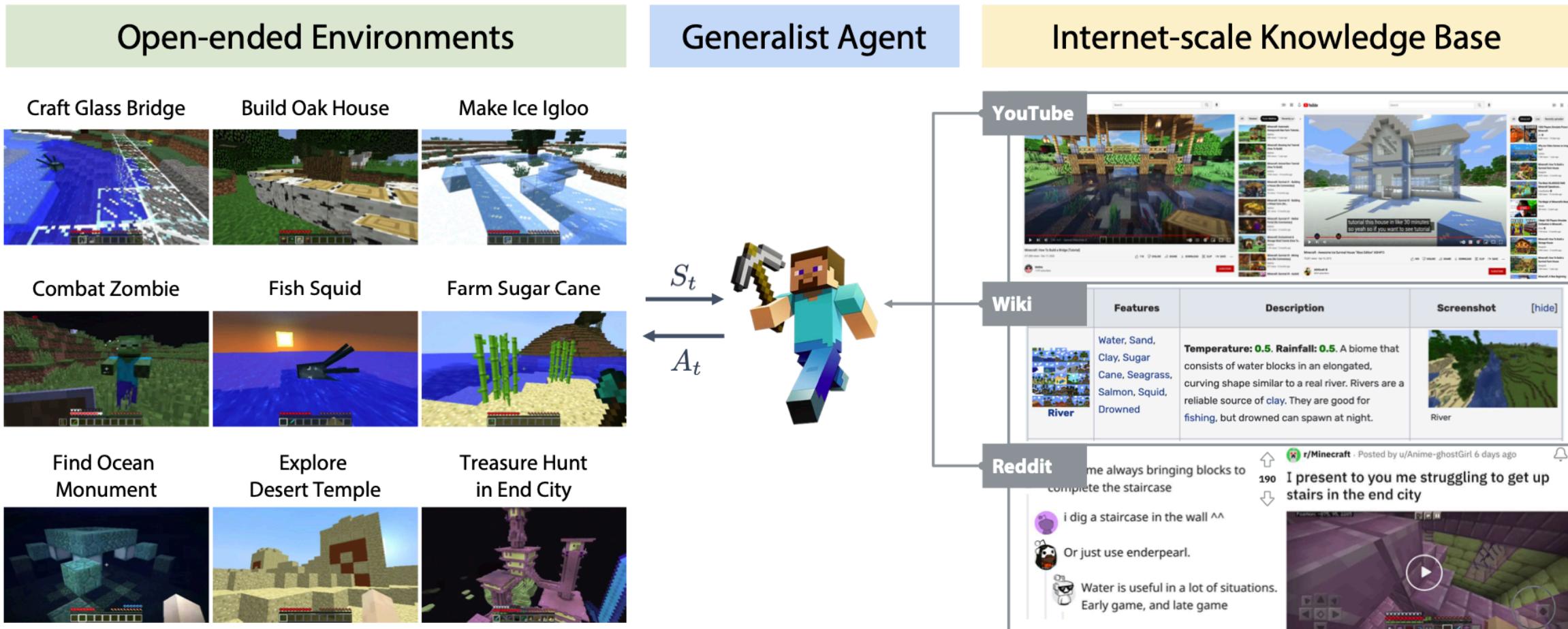














Heading OCR

Question: Tell me the heading text of this screenshot of webpage.

Answer: Discover, Appreciate, & Understand the Animal World!

Captioning

Question: What is the meta description of this website?

Answer: The world's largest & most trusted collection of animal facts, pictures and more!

WebQA

Question: What additional platform is mentioned for following the website's content?

Answer: YouTube Channel

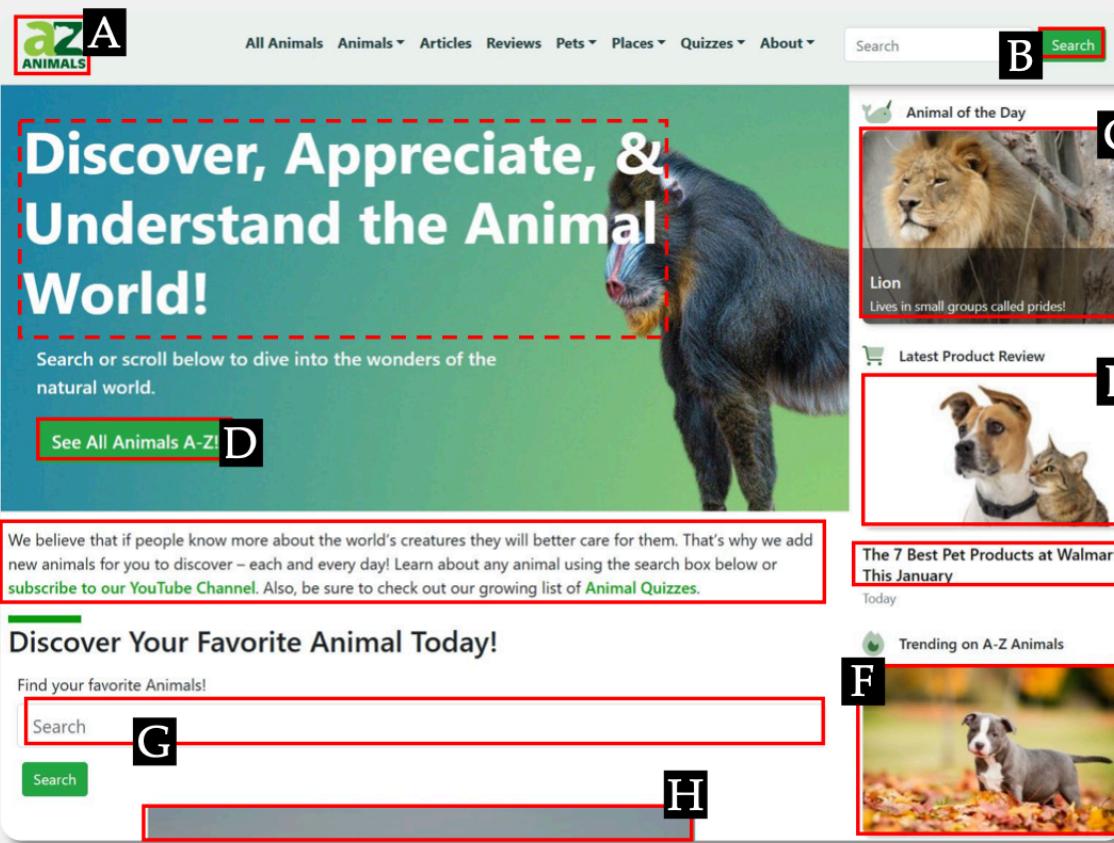
VisualWebBench

- Website-wise Task
- Element-wise Task
- Action-wise Task

Element OCR

Question: Tell me the text content in the red bounding box

Answer: We believe that if people know about the world's creatures they will better care for them. That's why we add new animals for you to discover ...



Element Grounding

Question: I have labeled bright IDs for some HTML elements in this website screenshot. Tell me which one is the element corresponding to the description: button with text "See All Animals A-Z!"

Answer: D

Action Grounding

Question: I have labeled bright IDs for some HTML elements in this website screenshot. Tell me which one I should click to complete the instruction: learn about the animal of the day

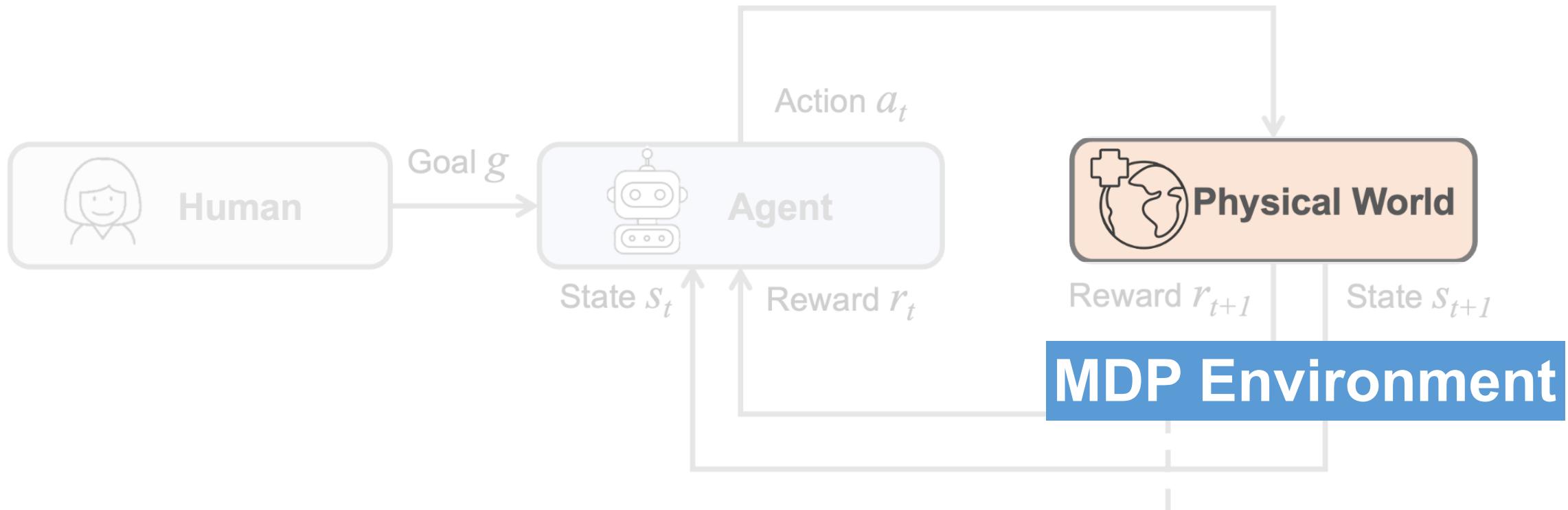
Answer: C

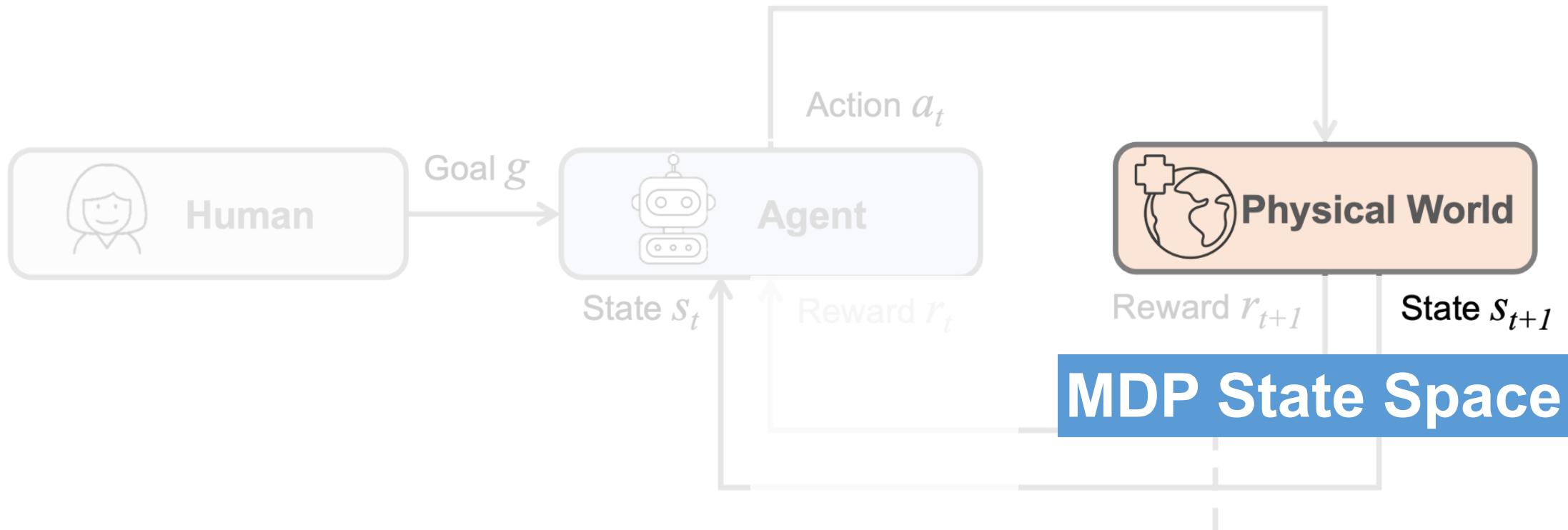
Action Prediction

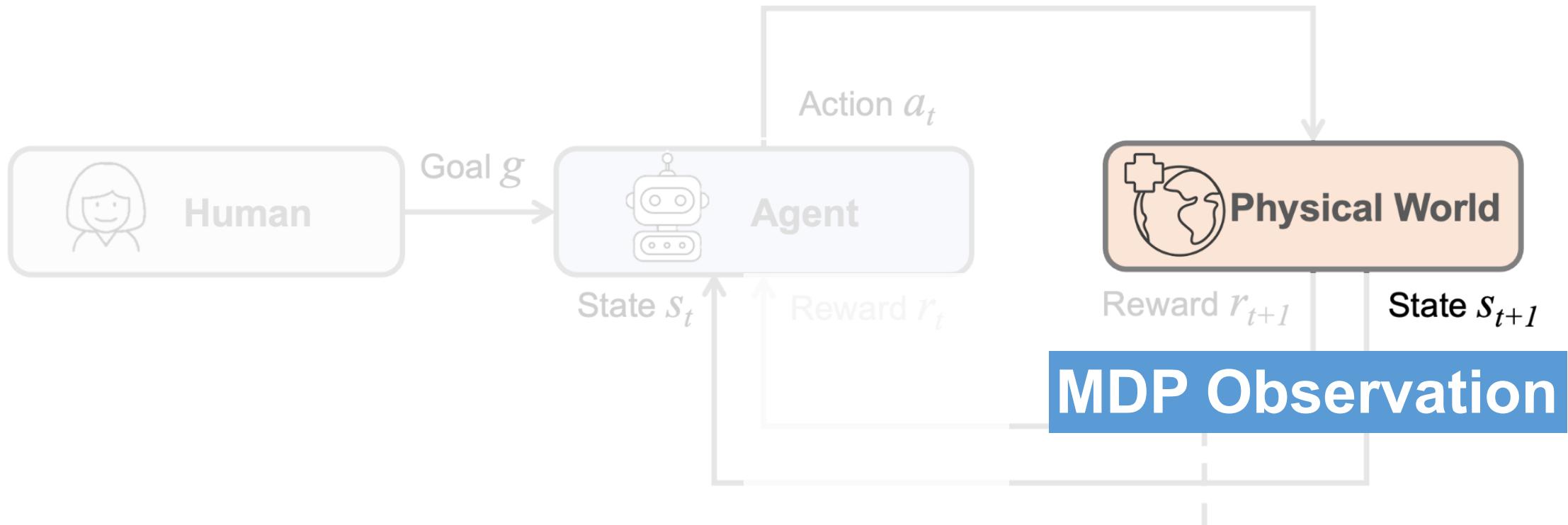
Question: After clicking the element in the bounding box, which one is the best description of the new webpage?

- Animal news, facts, ...
- All animals A-Z List
- The 7 best pet ...
- Search any animals!

Answer: C









Environment : Observation



Environment : Observation (Rendered 2D Images)



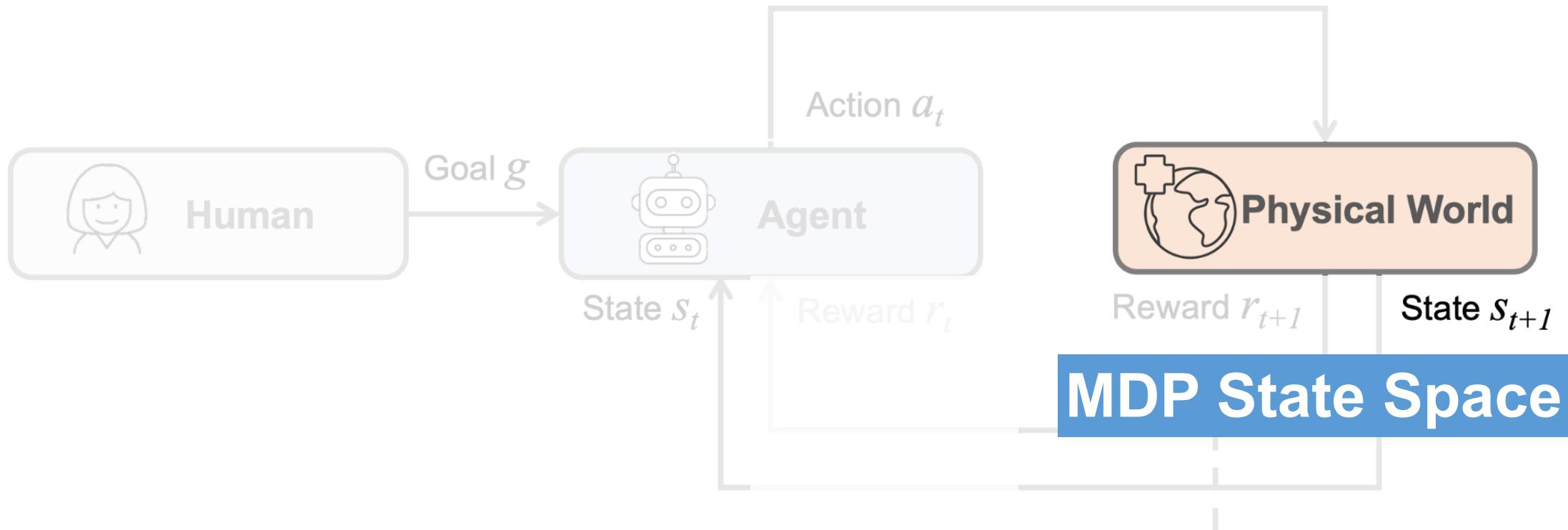
Environment : Observation (Rendered 2D Images)

Enabled by large dataset of realistic interactive **scenes** and objects



50 Scenes







Environment : Observation → State (3D Assets & States)



Environment : Observation → State (3D Assets & States)

50 Scenes

Enabled by large dataset of realistic interactive scenes and objects



10000 Objects

Semantic

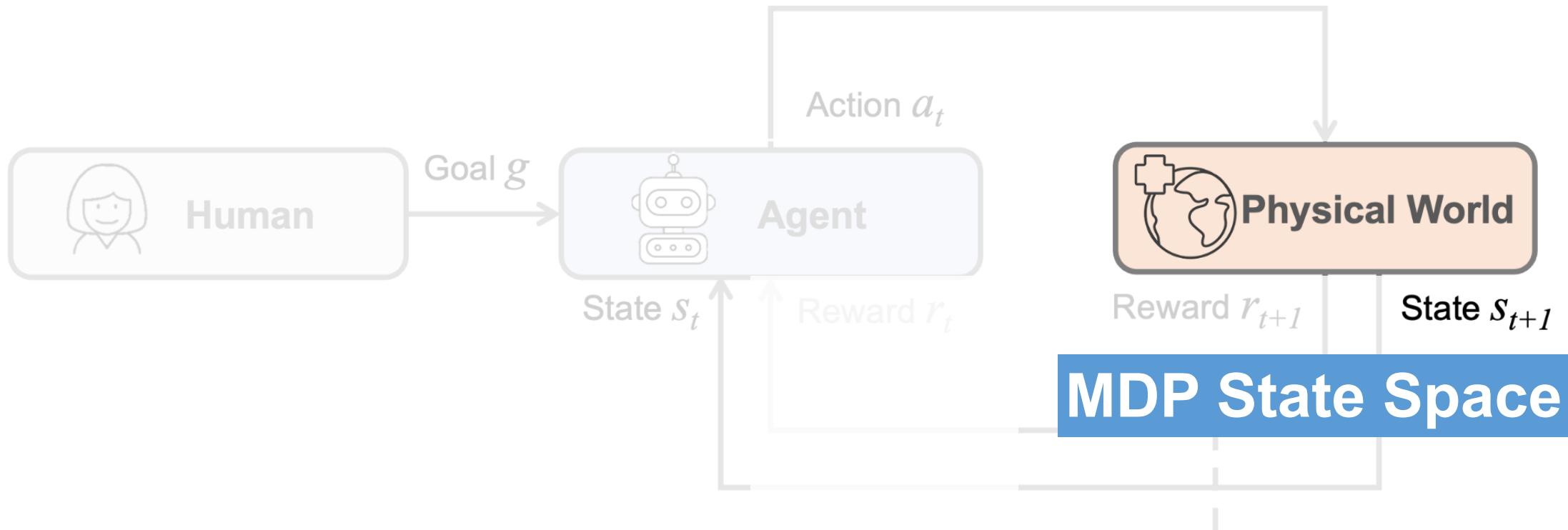
Properties: cookable, sliceable, freezable, burnable, deformable

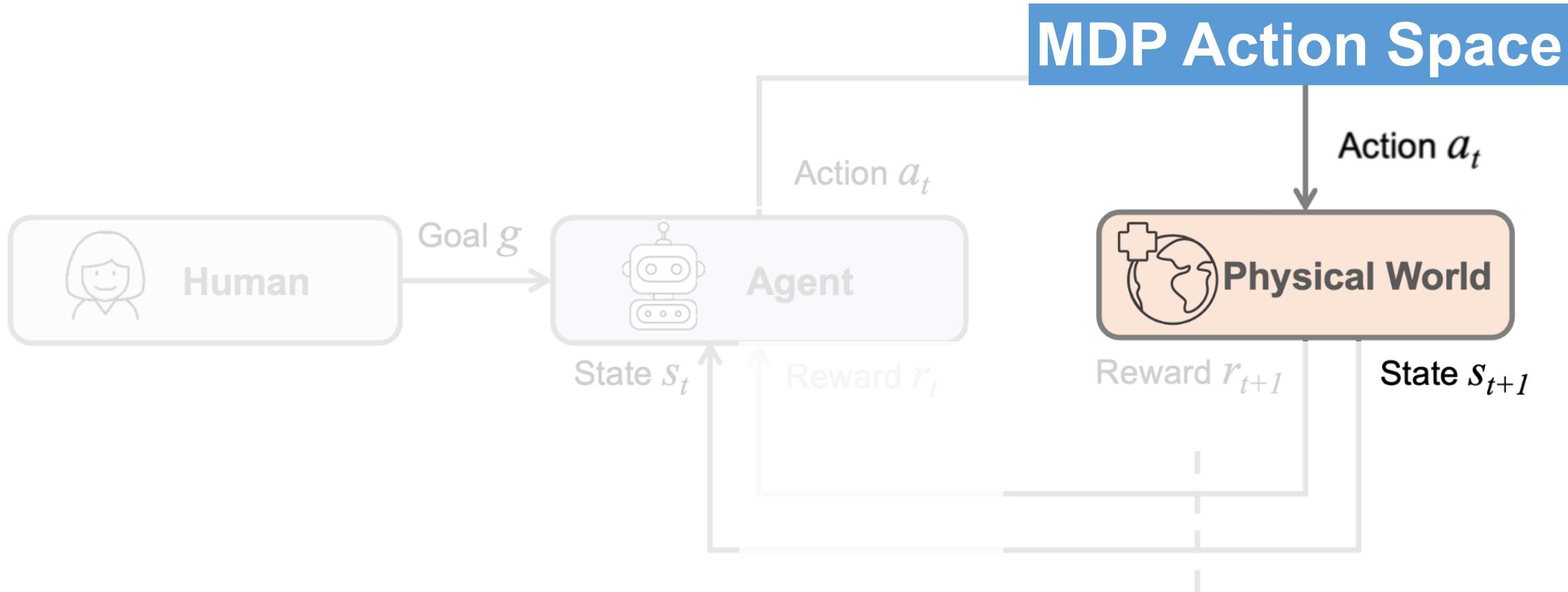
...
Cooking temperature: 58°C

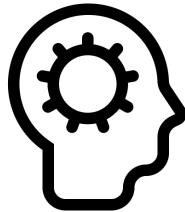
Physical



Articulation annotation
(joint type, origin, axis, limit)
Mass, friction, CoM, ...
Canonical size and orientation

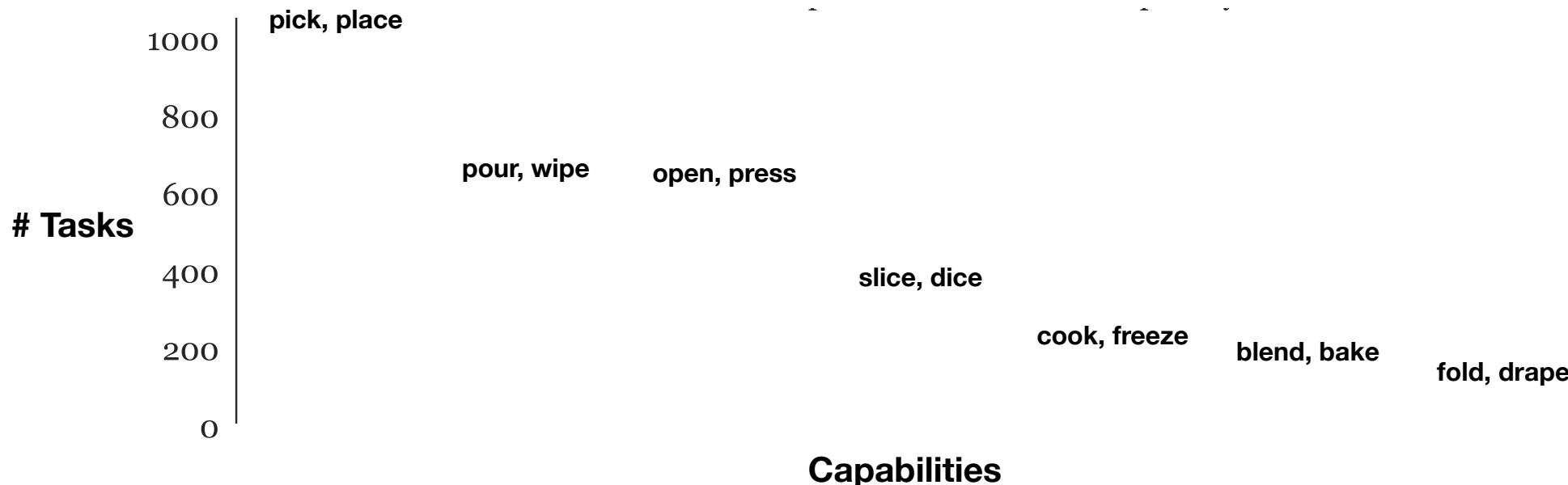


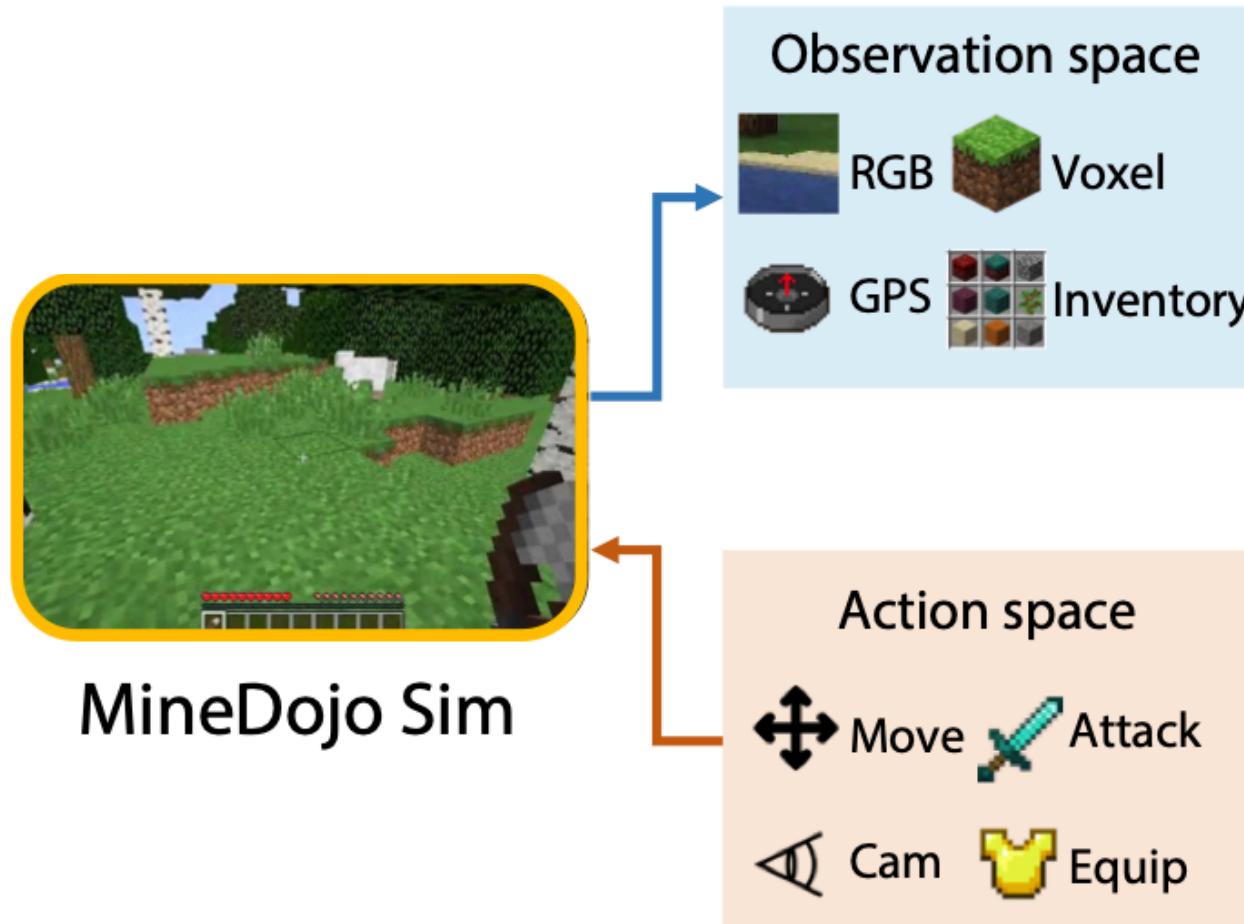




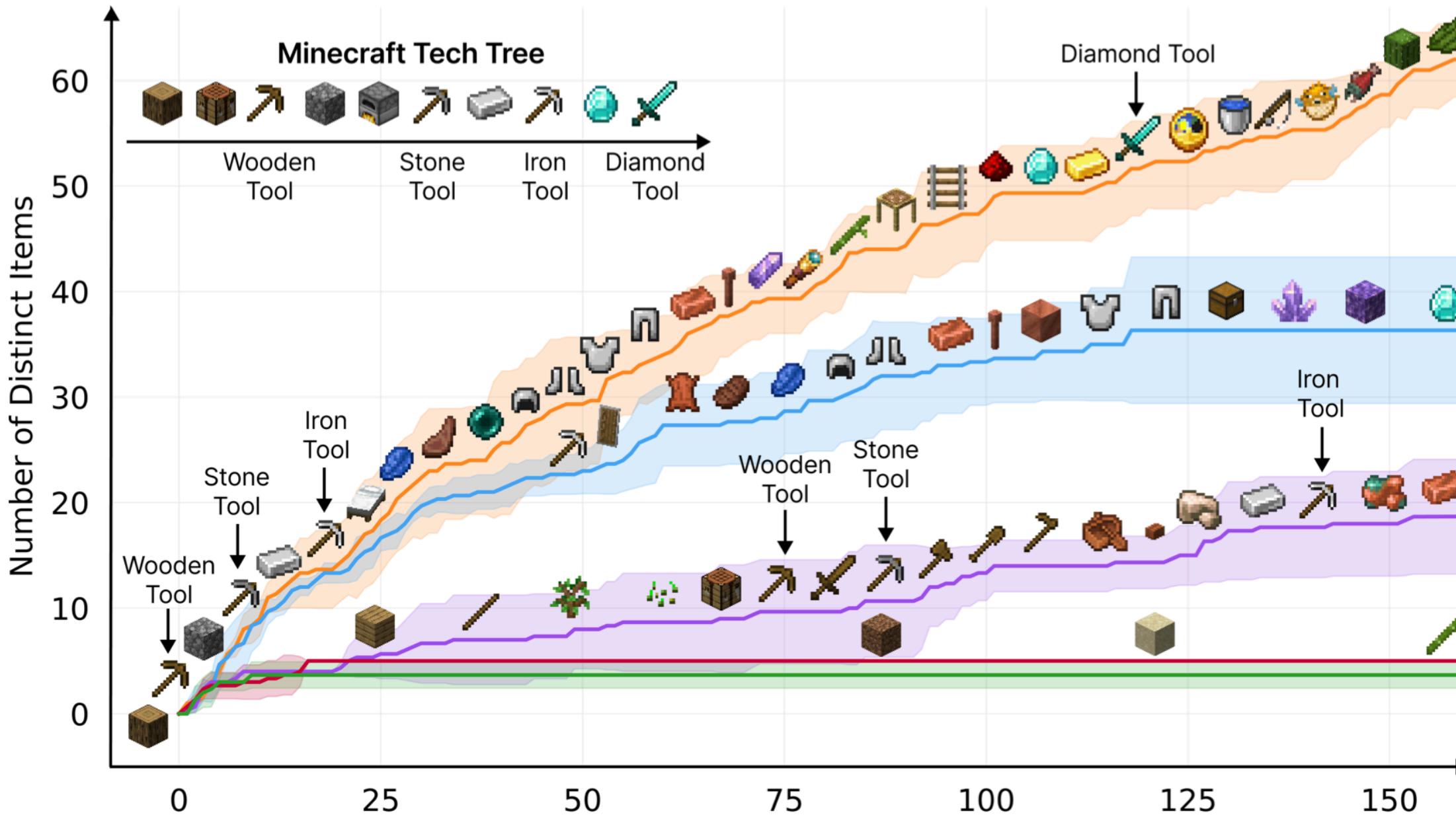
Action : can robots learn to **solve** these tasks?

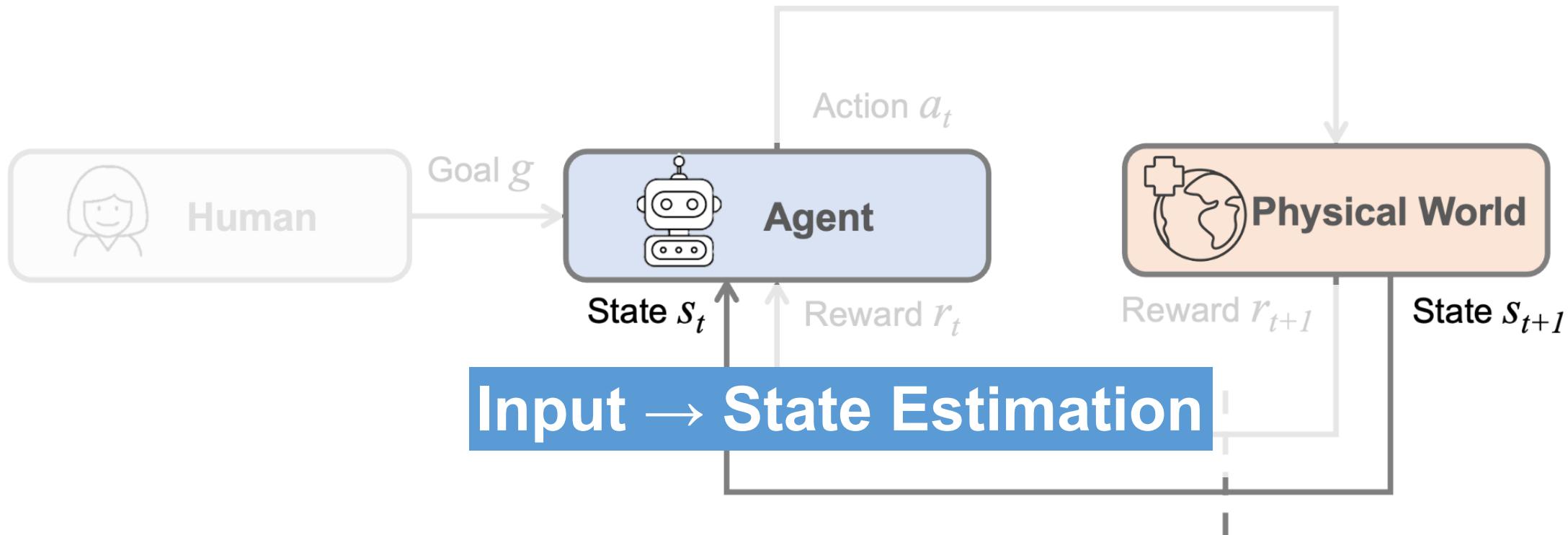
→ What **capabilities** are needed?





MDP Action Space: Skills





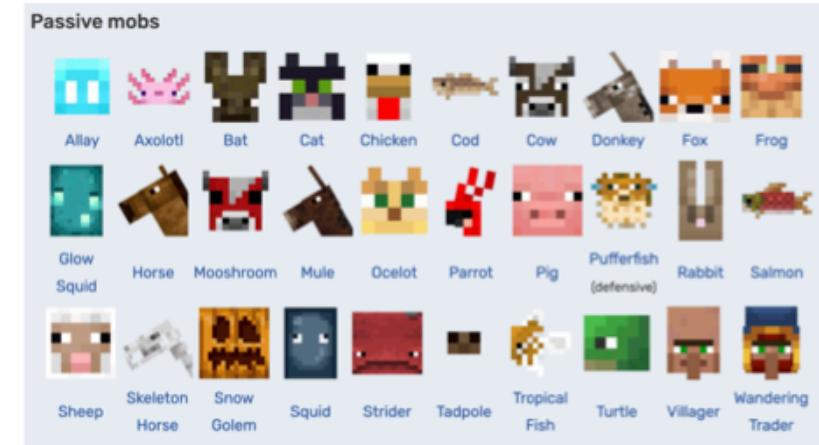
Perception / State Estimation

$$o \rightarrow s$$

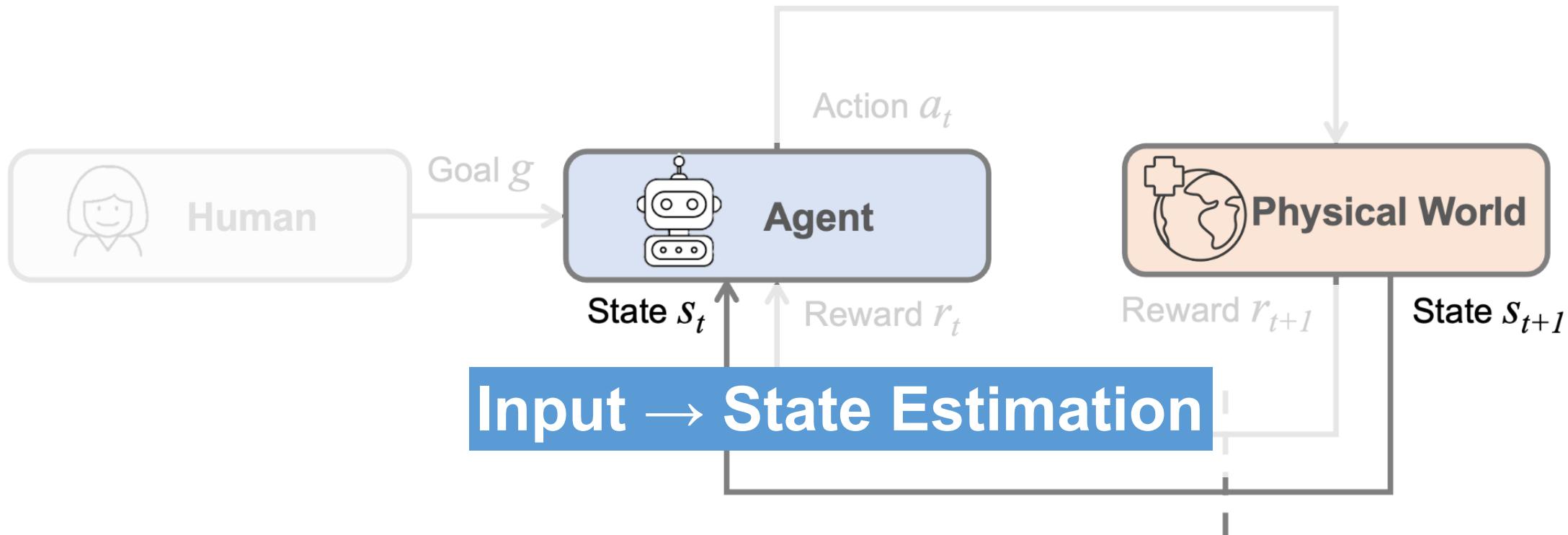
Observation (2D rendered scenes)

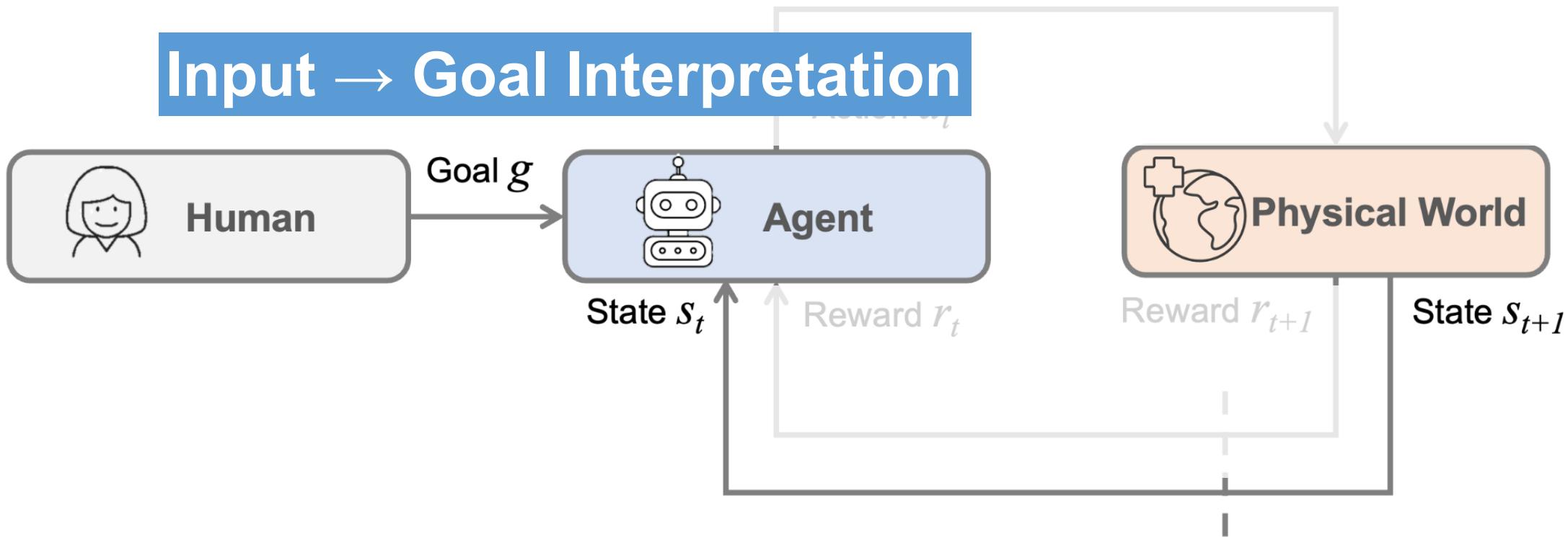


State (3D assets)



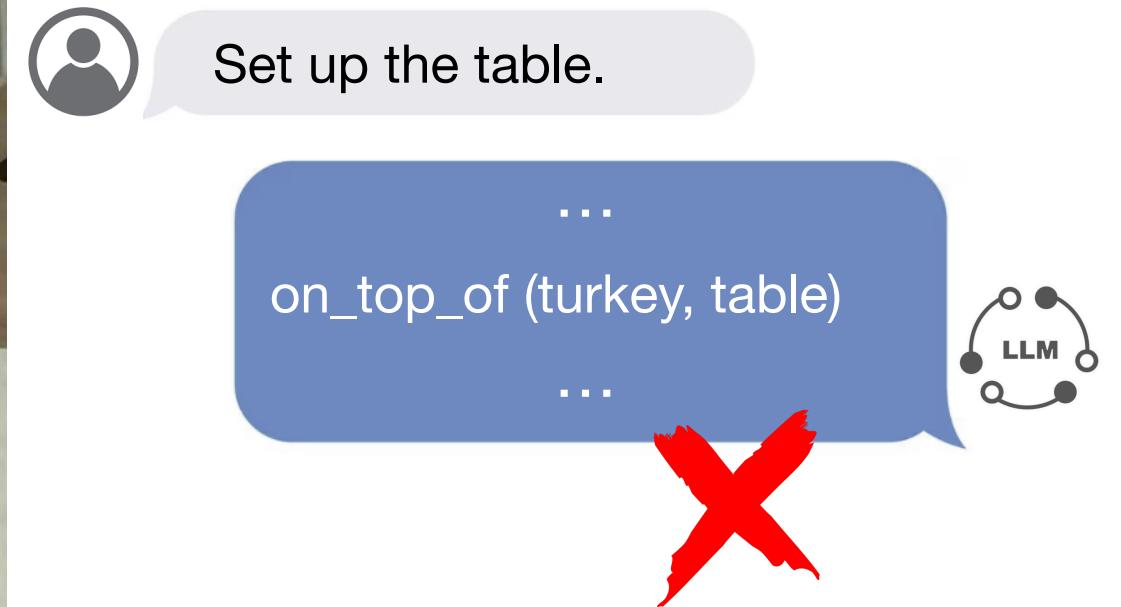
Name	Ingredients	Crafting recipe
Cake	Milk Bucket + Sugar + Egg + Wheat	
Golden Apple	Gold Ingot + Apple	





Goal Interpretation

g





Set up the table.

...

on_top_of (turkey, table)

...



Use the plates.

...

on_top_of (plate, table)

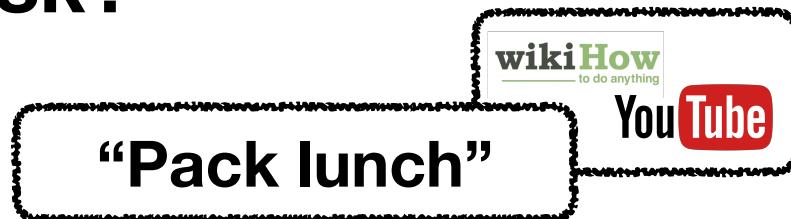
...

on_top_of (turkey, plate)





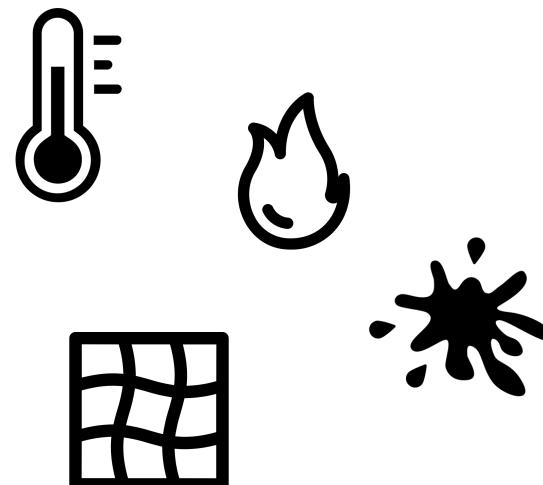
Goal : defines a task?



What objects?



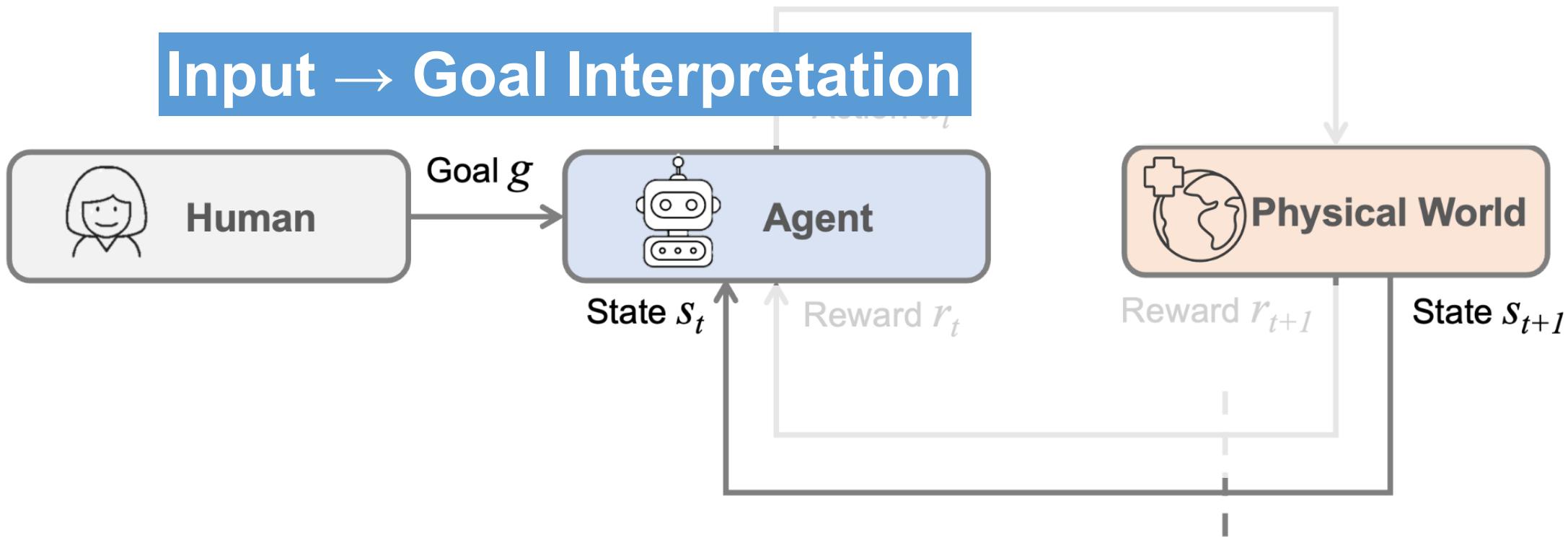
What properties?

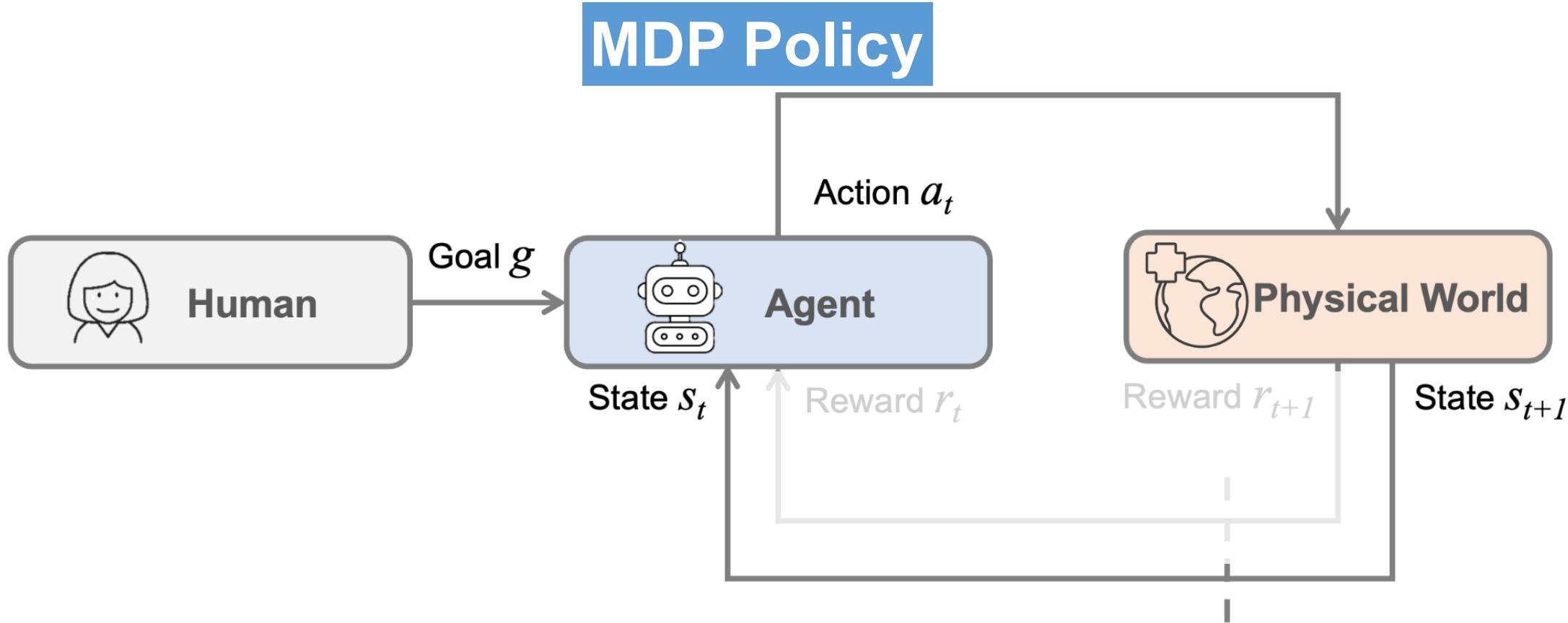


Start & Goal?

- apple in fridge
- burger in fridge
- water bottle in fridge
- paper bag on counter







Policy

$$\pi(o, g) \rightarrow a$$

Input: Preserving food

Environment

... inside (strawberry, pan) ...



Goal

... cooked (strawberry) ...



Operator

BEHAVIOR**Action Trajectory**

- ... ✓
- A7 OPEN(oven) ✓
- A8 RIGHT_GRASP(pan) ✓
- A9 RIGHT_PLACE_INSIDE(oven)
- A10 CLOSE(oven)
- A11 COOK(strawberry)
- ...

LLM Output

This video is for demonstration only. There're no actual controller-level actions. For action execution examples, visit our repository: <https://github.com/embodied-agent-interface/embodied-agent-interface>.

Input: Preserving food

Environment

... inside (strawberry, pan) ...



Goal

... cooked (strawberry) ...



Operator

BEHAVIOR**Action Trajectory**

- ... ✓
- A7 OPEN (oven) ✓
- A8 RIGHT_GRASP(pan) ✓
- A9 RIGHT_PLACE_INSIDE (oven) ✓
- A10 CLOSE(oven)
- A11 COOK(strawberry)
- ...

LLM Output

This video is for demonstration only. There're no actual controller-level actions. For action execution examples, visit our repository: <https://github.com/embodied-agent-interface/embodied-agent-interface>.

Input: Preserving food

Environment

... inside (strawberry, pan) ...



Goal

... cooked (strawberry) ...



Operator

BEHAVIOR**Action Trajectory**

- ... ✓
- A7 OPEN (oven) ✓
- A8 RIGHT_GRASP(pan) ✓
- A9 RIGHT_PLACE_INSIDE (oven) ✓
- A10 CLOSE(oven) ✓
- A11 COOK(strawberry)
- ...

LLM Output

This video is for demonstration only. There're no actual controller-level actions. For action execution examples, visit our repository: <https://github.com/embodied-agent-interface/embodied-agent-interface>.

Input: Preserving food

Environment

... inside (strawberry, pan) ...



Goal

... cooked (strawberry) ...



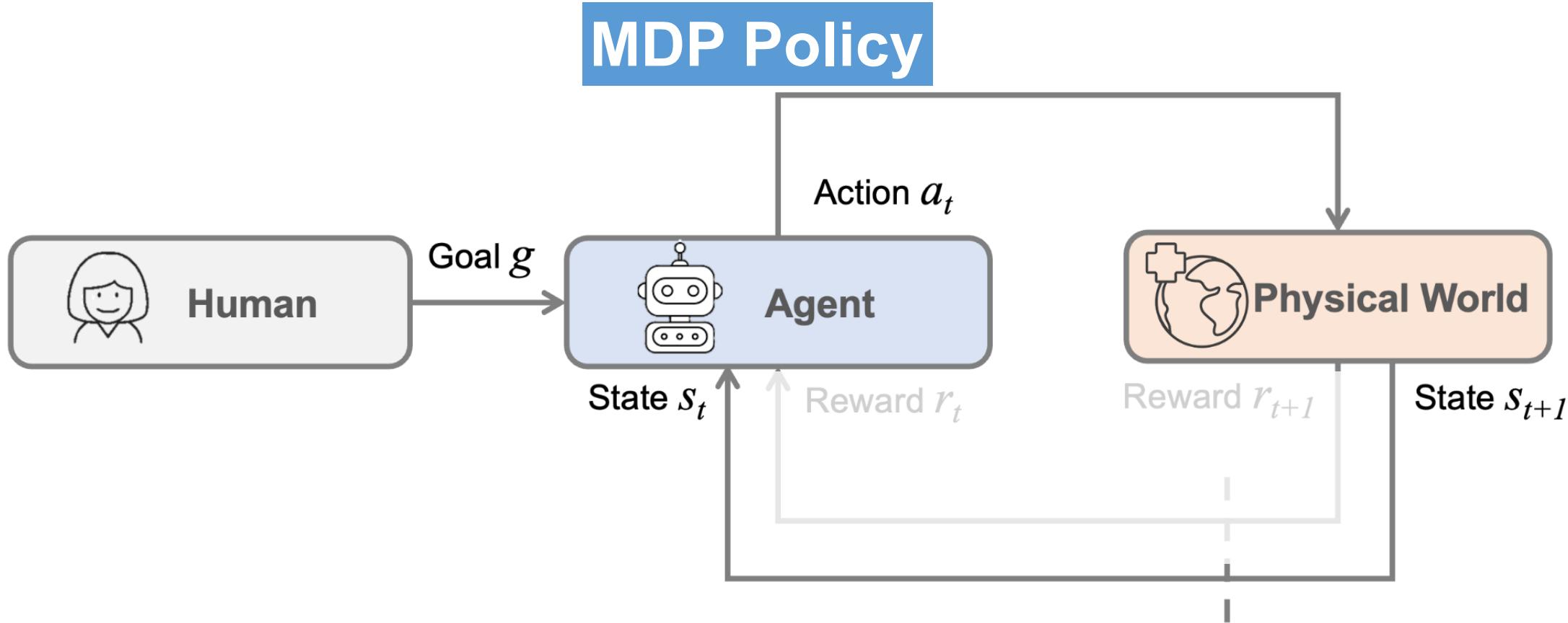
Operator

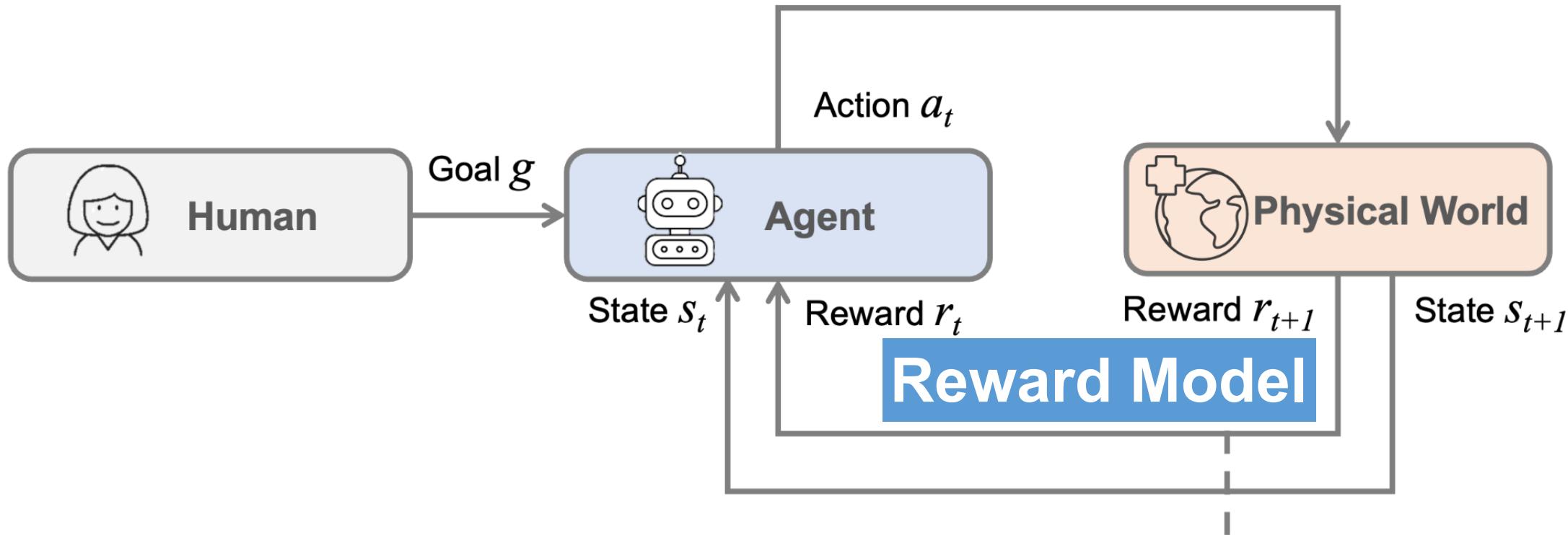
BEHAVIOR**Action Trajectory**

- ... ✓
- A7 OPEN (oven) ✓
- A8 RIGHT_GRASP(pan) ✓
- A9 RIGHT_PLACE_INSIDE (oven) ✓
- A10 CLOSE(oven) ✓
- A11 COOK(strawberry) ✓
- ...

LLM Output

This video is for demonstration only. There're no actual controller-level actions. For action execution examples, visit our repository: <https://github.com/embodied-agent-interface/embodied-agent-interface>.



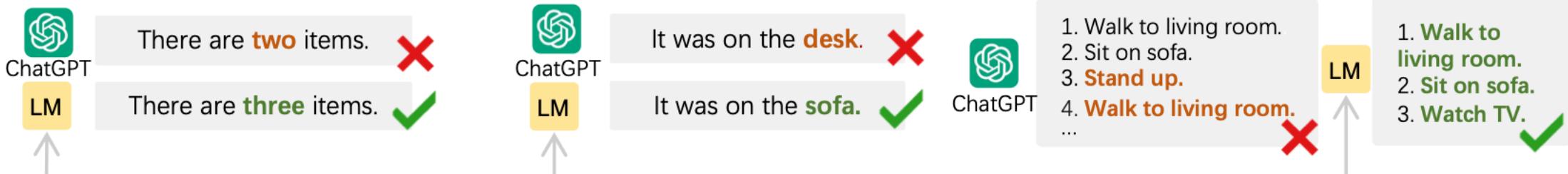


Reward Model

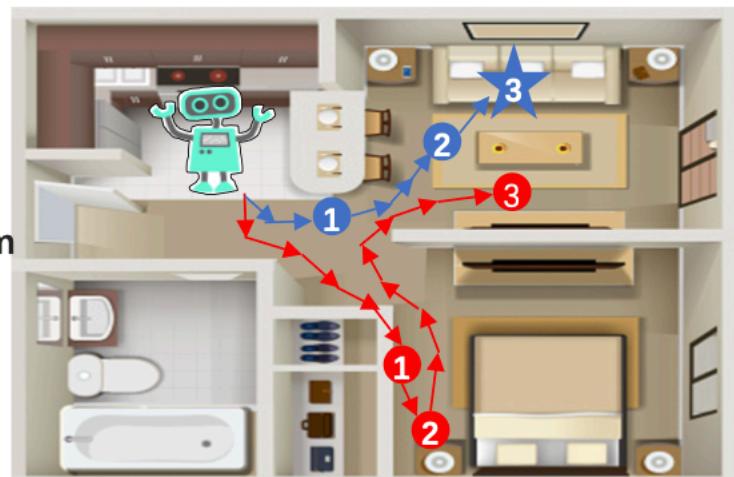
$$o, a \rightarrow r$$

Reward Model

Answer:



E2WM
training paradigm



World Model

Goal-Oriented
Planning

Random
Exploration

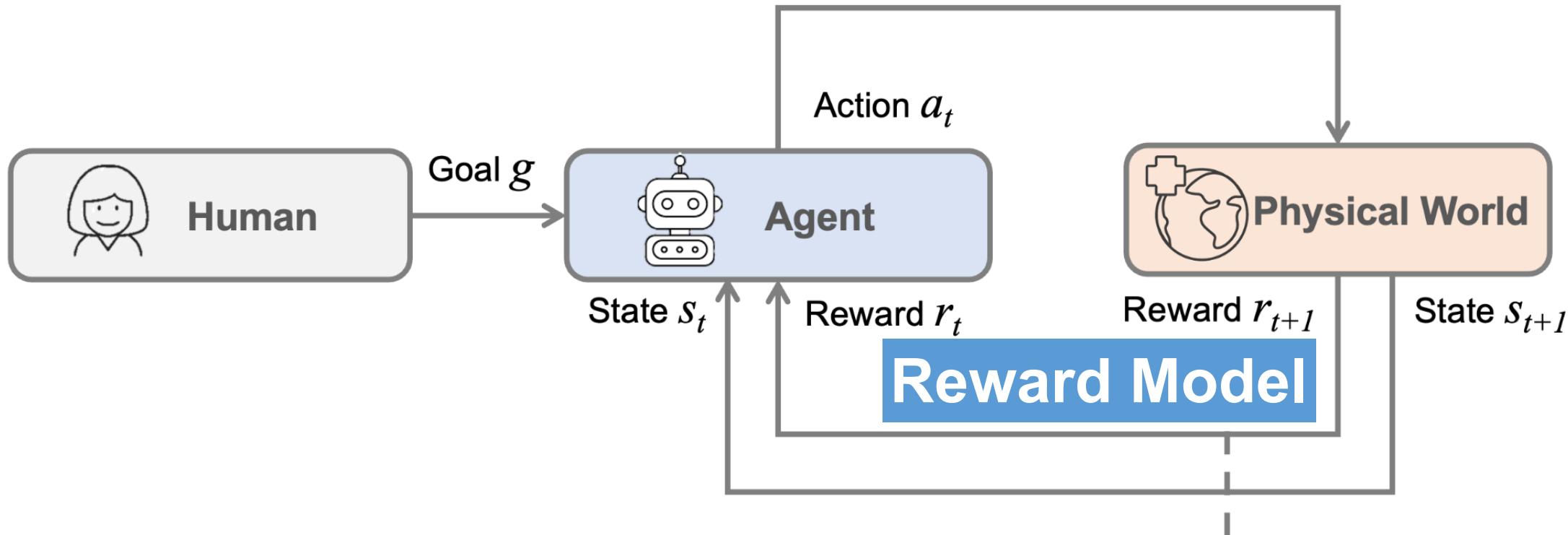
Embodied Experiences

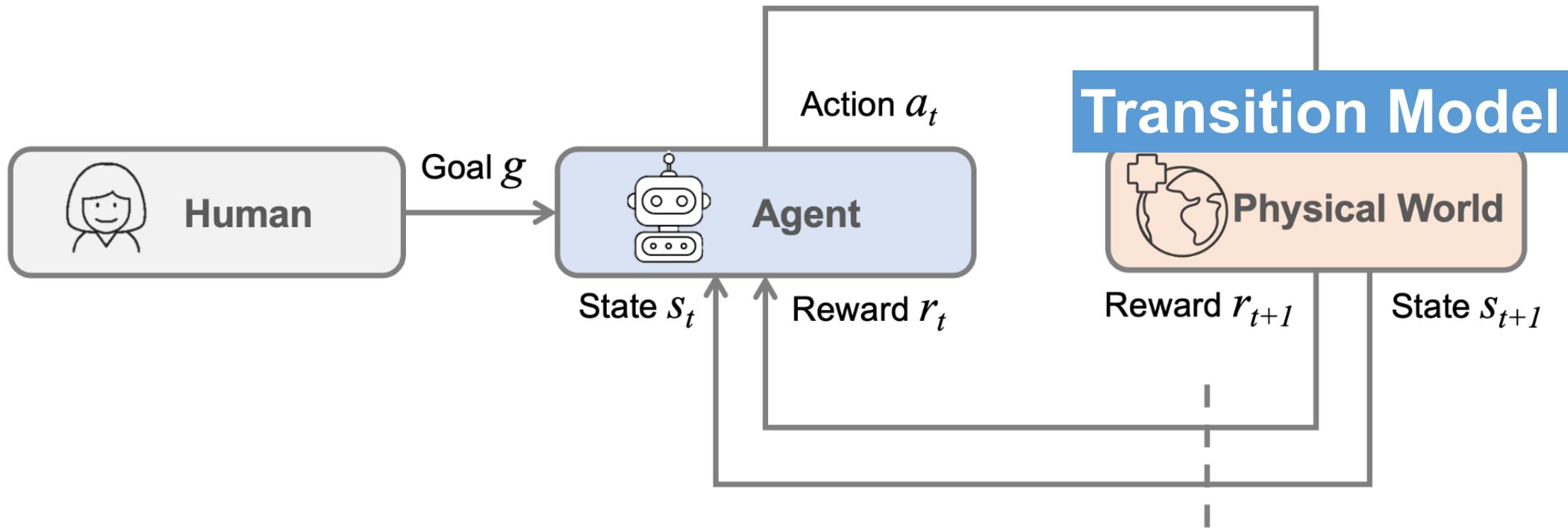
Activity: Watch TV
① Walk to living room.
② Sit on sofa.
③ Watch TV.

Action
① Walk to bedroom.
② Grab book.
③ Walk to living room.
Tracking (of book)
On bed, in bedroom
In hand, in bedroom
In hand, in living room.

Finetune

Language Model
(e.g. GPT-J-6B)





Transition Model

$$o_t, a \rightarrow o_{t+1}$$

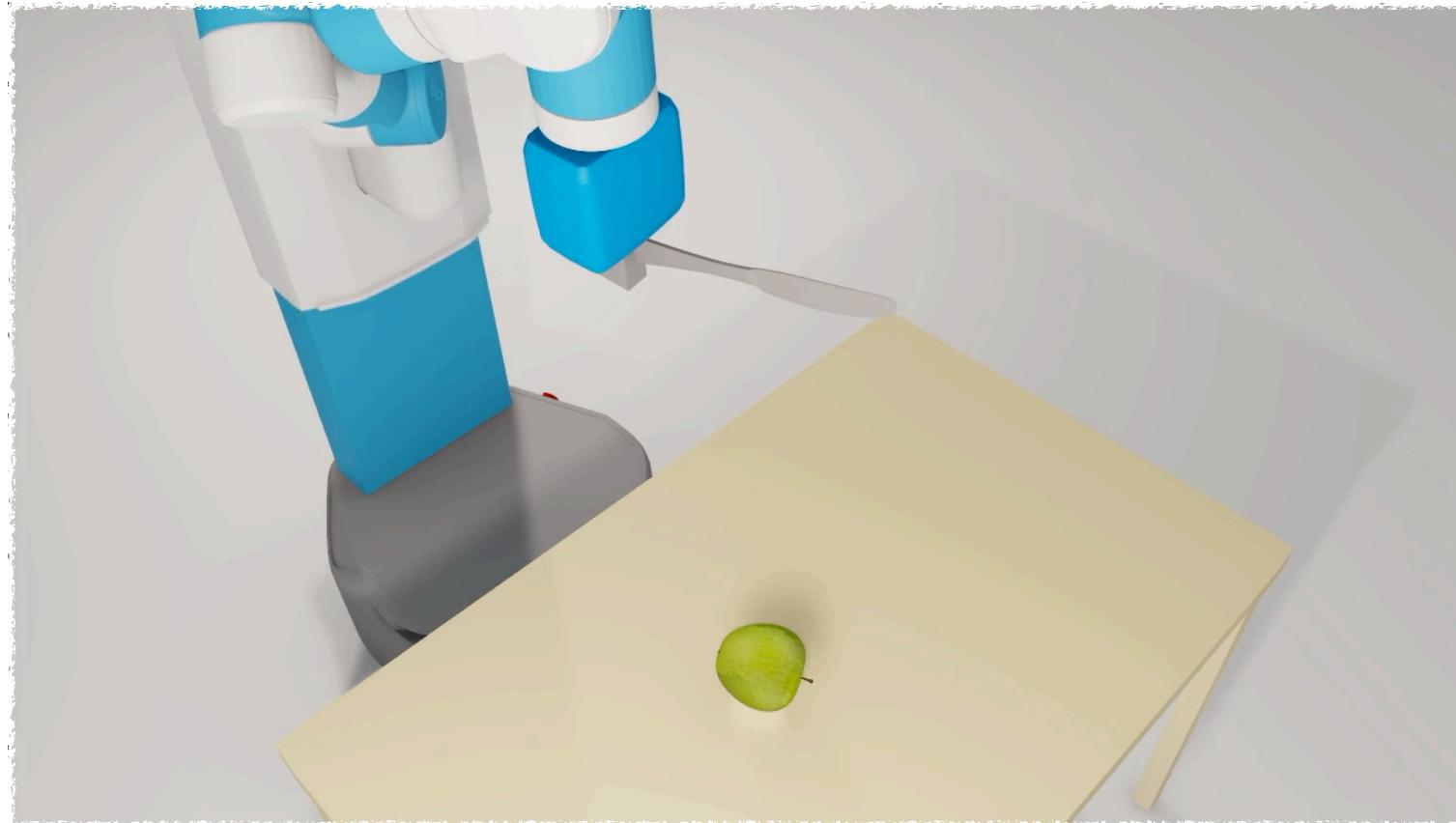
"World Modeling"



OmniGibson

Physics Transition Examples

```
class SlicingRule:
```





OmniGibson

Physics Transition Interface

```
class TransitionRule:  
    def condition(self, *args) -> bool  
    def transition(self, *args)
```

Determines whether a transition should occur

What should happen when a transition is triggered

Allows us to capture arbitrarily complex **physical phenomena!**



OmniGibson

Physics Transition Examples

```
class SlicingRule:

    def condition(self, *args) -> bool
        # Return True if a slicer object is touching
        # a sliceable object with sufficient force

    def transition(self, *args)
        # Remove the sliceable object and
        # import its sliced component objects
```

1

Large-Scale
Scene Generation

2

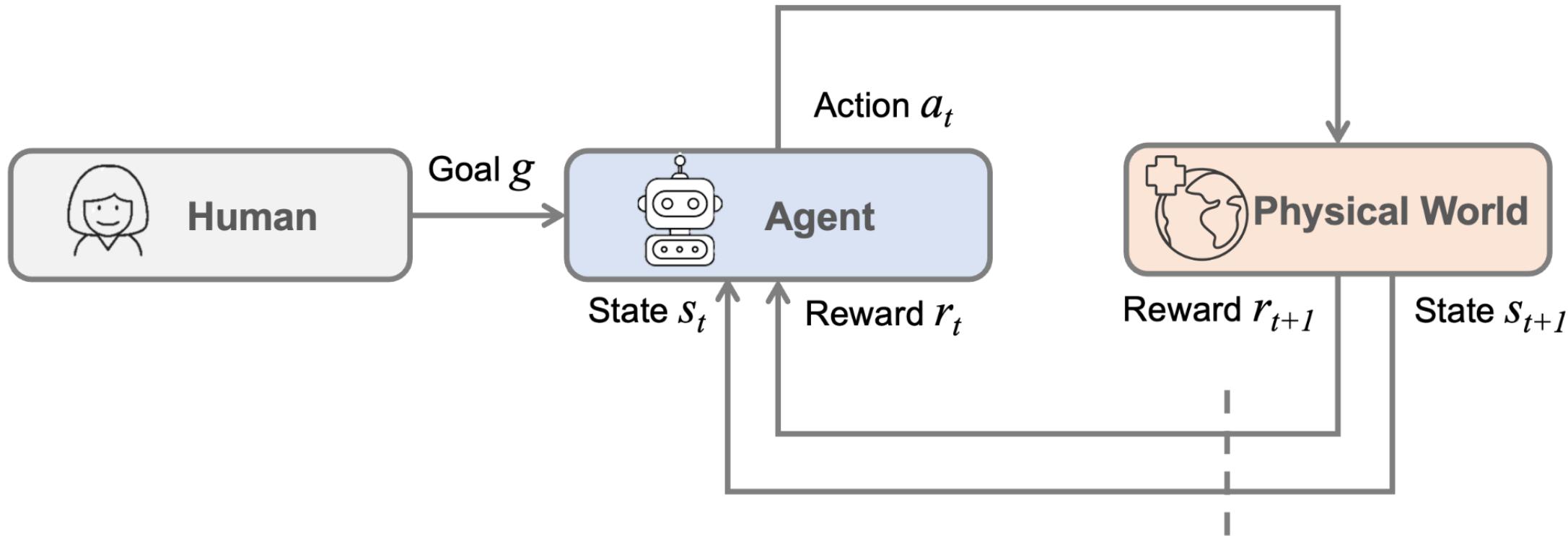
Modular
Robot Control

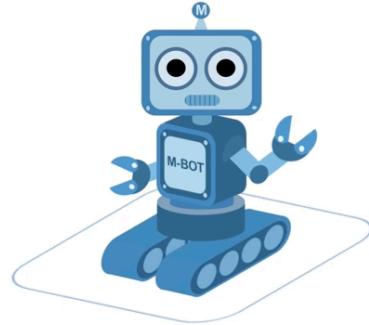
3

Kinematic & Semantic
Object States

4

Customizable
Physics Transitions





Different Instantiations of MDP



Perfect actuations

Noisy actuations

Perfect odometry

Noisy odometry

RGB only perception

Multiple perceptual modalities

-
-
-

-
-
-

Outline



Content	Time	Presenter
1. Motivation and Overview	15min	Manling Li
2. Foundation Models meet Virtual Agents	45min	Manling Li
3. Foundation Models meet Physical Agents Overview & Perception High-level & Low-level Decision Making	25min 50min	Jiayuan Mao Wenlong Huang
4. Robotic Foundation Models	30min	Yunzhu Li
5. Remaining Challenges	15min	Yunzhu Li
QA	30min	

Outline



Content	Time	Presenter
1. Motivation and Overview	15min	Manling Li
2. Foundation Models meet Virtual Agents	45min	Manling Li
3. Foundation Models meet Physical Agents Overview & Perception High-level & Low-level Decision Making	25min 50min	Jiayuan Mao Wenlong Huang
4. Robotic Foundation Models	30min	Yunzhu Li
5. Remaining Challenges	15min	Yunzhu Li
QA	30min	

Outline



Content	Time	Presenter
1. Motivation and Overview	15min	Manling Li
2. Foundation Models meet Virtual Agents	45min	Manling Li
3. Foundation Models meet Physical Agents Overview & Perception High-level & Low-level Decision Making	25min 50min	Jiayuan Mao Wenlong Huang
4. Robotic Foundation Models	30min	Yunzhu Li
5. Remaining Challenges	15min	Yunzhu Li
QA	30min	