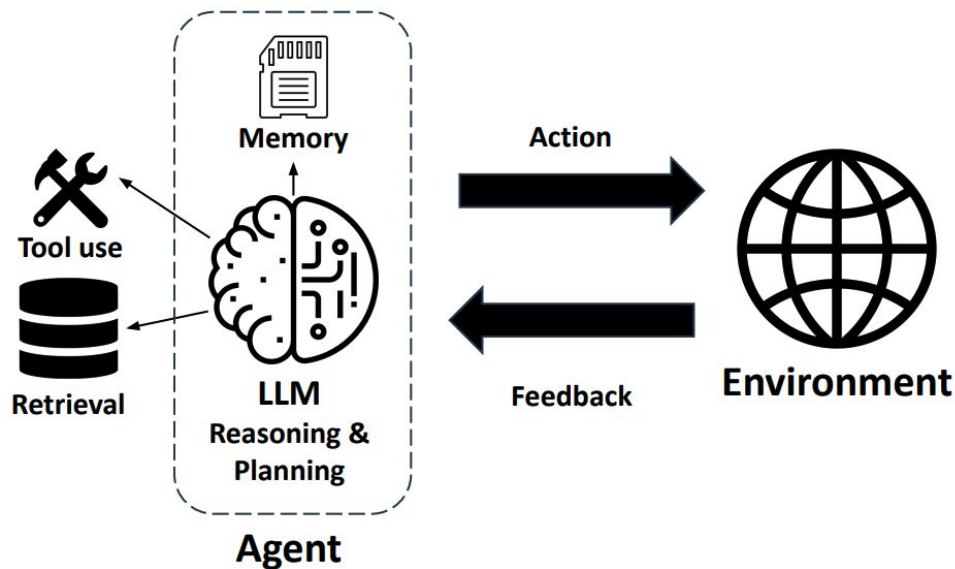# LLM Agent

Advanced Topics in Embodied Learning and Vision
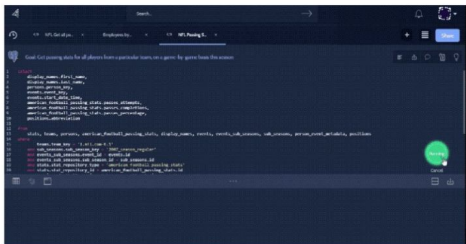
Ying Wang

2025.02.27

# LLM Agents

"An agent is anything that can be viewed as perceiving its environment through sensors and acting upon that environment through actuators." — Russell & Norvig, AI: A Modern Approach (2020)
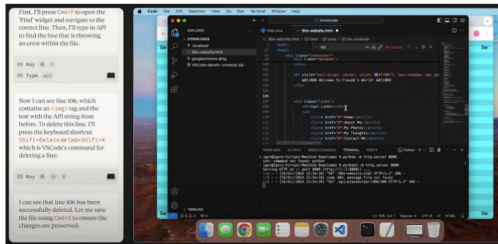
# **Why LLM Agents?**

➔  Solving real-world tasks typically involves a trial-and-error process

➔  Leveraging external tools and retrieving from external knowledge expand LLM's capabilities

➔  Agent workflow facilitates complex tasks

- Task decomposition
- Allocation of subtasks to specialized modules
- Division of labor for project collaboration
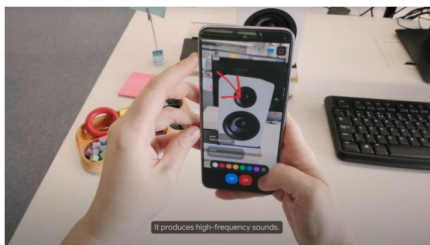- Multi-agent generation inspires better responses

3

# Applications

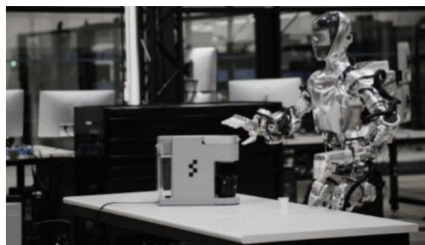**Code generation**
Cursor, GitHub Copilot, Devin, Google Jules…

**Computer use**
Anthropic Claude, Google Jarvis, OpenAI Operator

**Personal assistant**
Google Astra, OpenAI GPT-4o,…

**Robotics**
Figure AI, Tesla Optimus, NVIDIA GR00T…

- Education
- Law
- Finance
- Healthcare
- Cybersecurity
  …

*cr: https://rdi.berkeley.edu/adv-llm-agents/sp25*
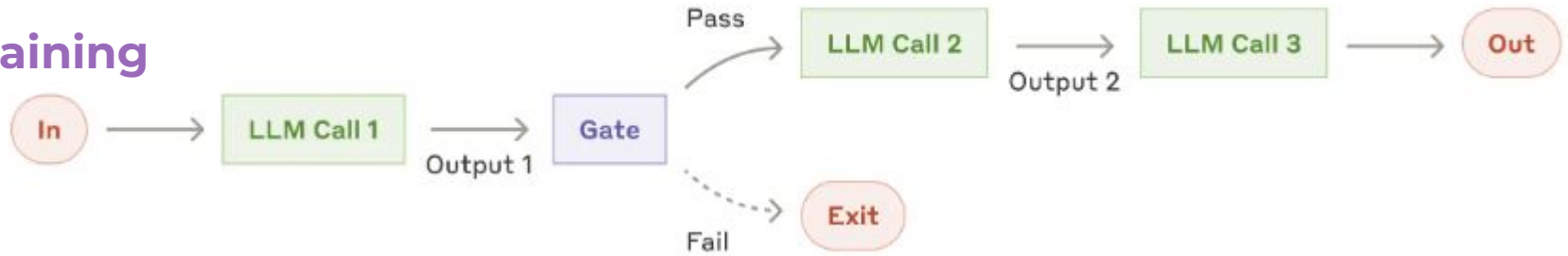
# Agenda

- Introduction
- Building block, workflows, agent
- Cognitive architectures for language agents
- LLM Agent Environments
    - Embodied Agent Interface
    - AgentBoard

# Building Block

cr: https://www.anthropic.com/research/building-effective-agents

# Workflows

**chaining**



**routing**

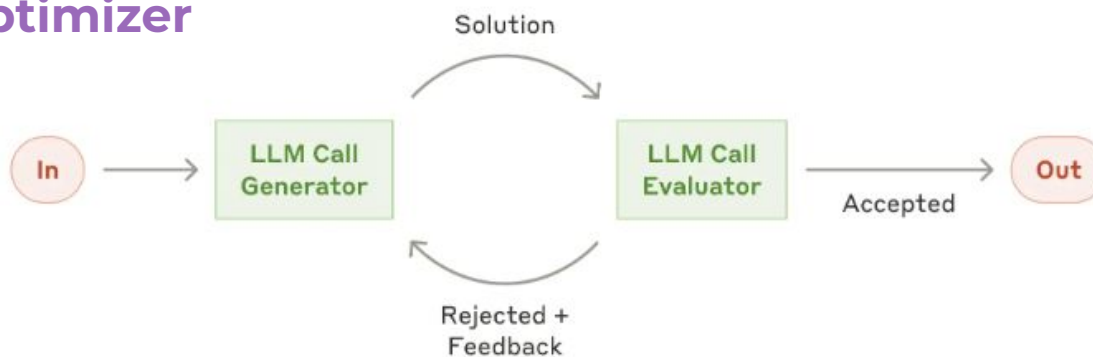*cr: https://www.anthropic.com/research/building-effective-agents*
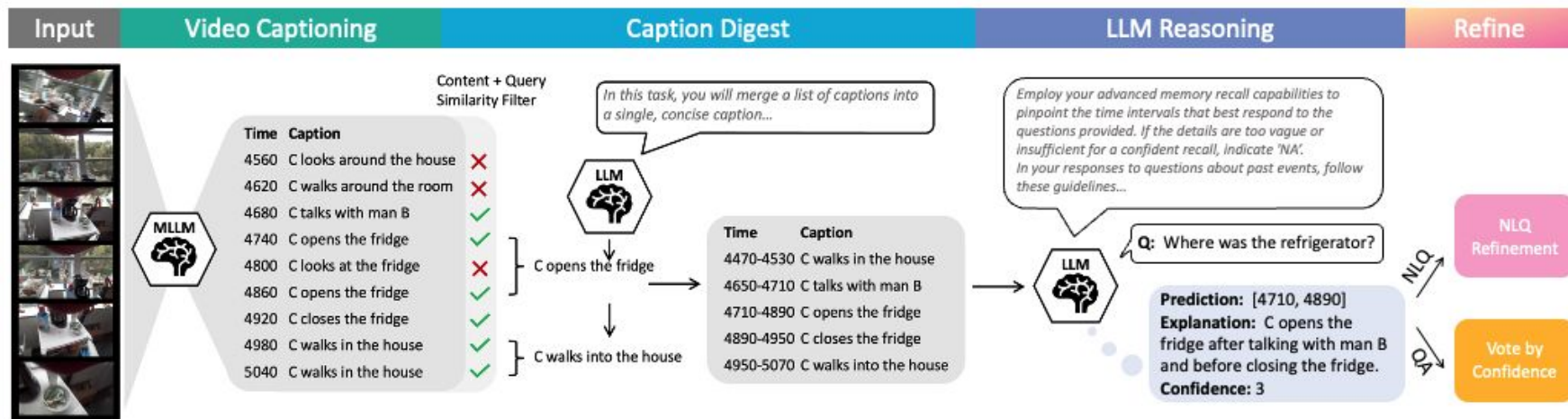
# Parallelization



Voting, sectioning…

# Evaluator-optimizer

# Example: LifelongMemory

- You may need to combine different workflows to solve a challenging task!

*cr: https://agenticlearning.ai/lifelong-memory/*

# Cognitive architectures for language agents

*cr: Cognitive Architectures for Language Agents (Sumers et al.)*

# CoALA: Memory

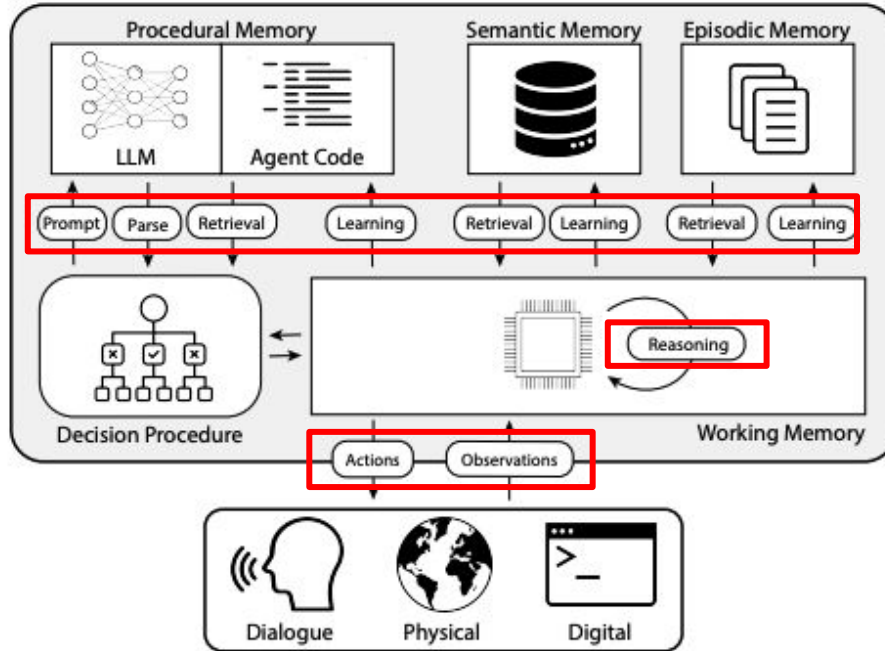

**Episodic memory** stores experience from earlier decision cycles.

**Semantic memory** stores an agent's knowledge about the world and itself.

**Procedural memory**: (i) implicit knowledge stored in the LLM weights; (ii) explicit knowledge written in the agent's code.

**Working memory** maintains active and readily available information as symbolic variables for the current decision cycle
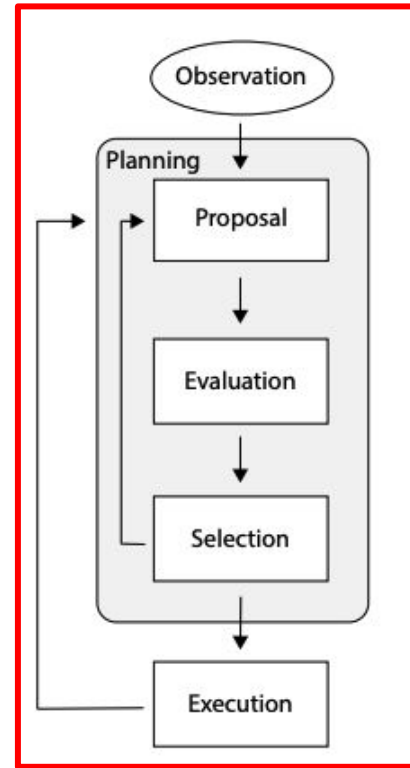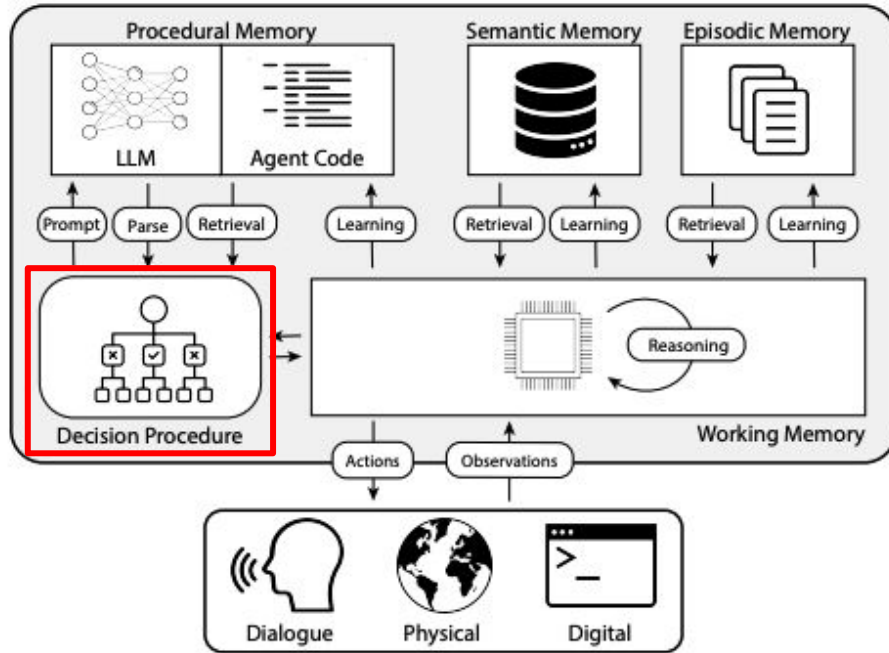
# CoALA: Action



**External actions** interact with external environments. E.g., control a robot, communicate with a human, navigate a website

**Internal actions** interact with internal memories.
- retrieval (read from long-term memory)
- learning (write to long-term memory)
- reasoning (update the short-term working memory with LLM)

# CoALA: Decision making

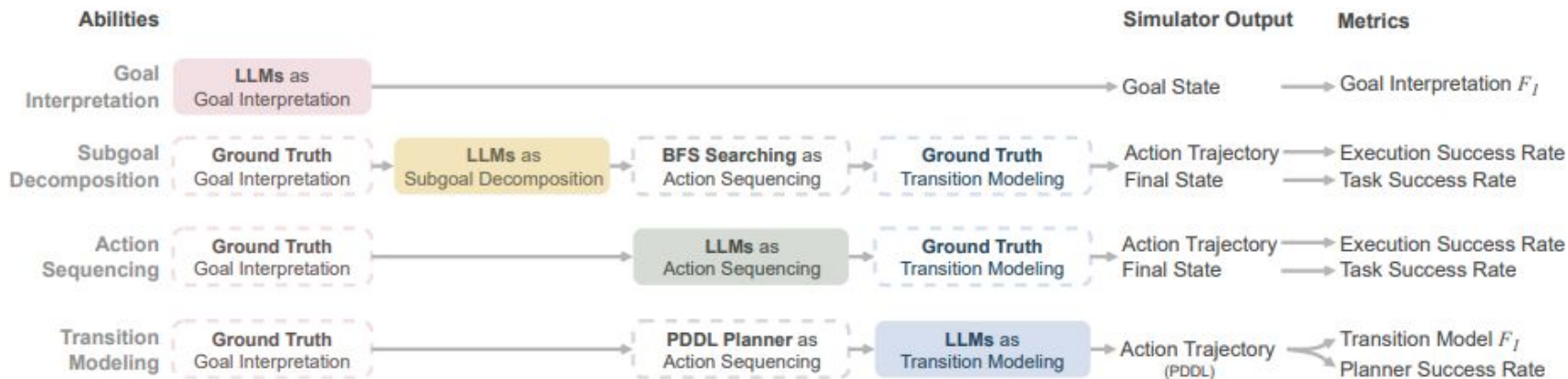# Agenda

- Introduction
- Building block, workflows, agent
- Cognitive architectures for language agents
- **LLM Agent Environments**
  - Embodied Agent Interface: Benchmarking LLMs for Embodied Decision Making (NeurIPS 2024) https://embodied-agent-interface.github.io/
  - AgentBoard: An Analytical Evaluation Board of Multi-turn LLM Agents (NeurIPS 2024) https://hkust-nlp.github.io/agentboard/

# Embodied Agent Interface: Benchmarking LLMs for Embodied Decision Making

Unlike existing evaluations rely solely on a final success rate, EAI breaks down the evaluation into

- **four modules** for decision making: goal interpretation, subgoal decomposition, action sequencing, and transition modeling
- a collection of **fine-grained metrics**, such as hallucination errors, affordance errors, various types of planning errors, etc.

# Embodied Agent Interface



For each ability module, to provide a comprehensive evaluation for it, we isolate this single module to be handled by the LLMs while using existing data or tools for the other modules.

# M1. Goal Interpretation

Ground the natural language instruction to the environment representations of objects, states, relations, and actions.



Goal Interpretation | Subgoal Decomposition | Action Sequencing | Transition Modeling

$g$

LTL Goal

Goal Interpretation

$\langle s_0, g_{nl} \rangle$

Initial State | Goal (Natural Language)

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

# M2. Subgoal Decomposition

Subgoal Decomposition generates a sequence of states, where each state can be a set of objects and their states.

# M3. Action Sequencing

Action Sequences are essential to achieve the state transitions identified in Subgoal Decomposition.
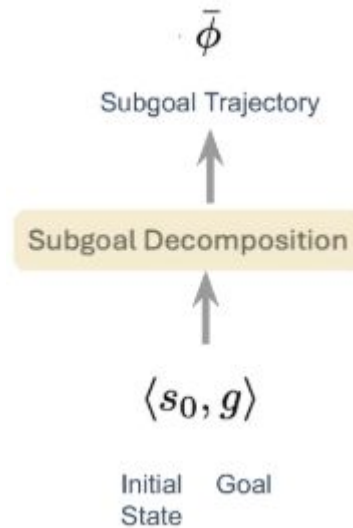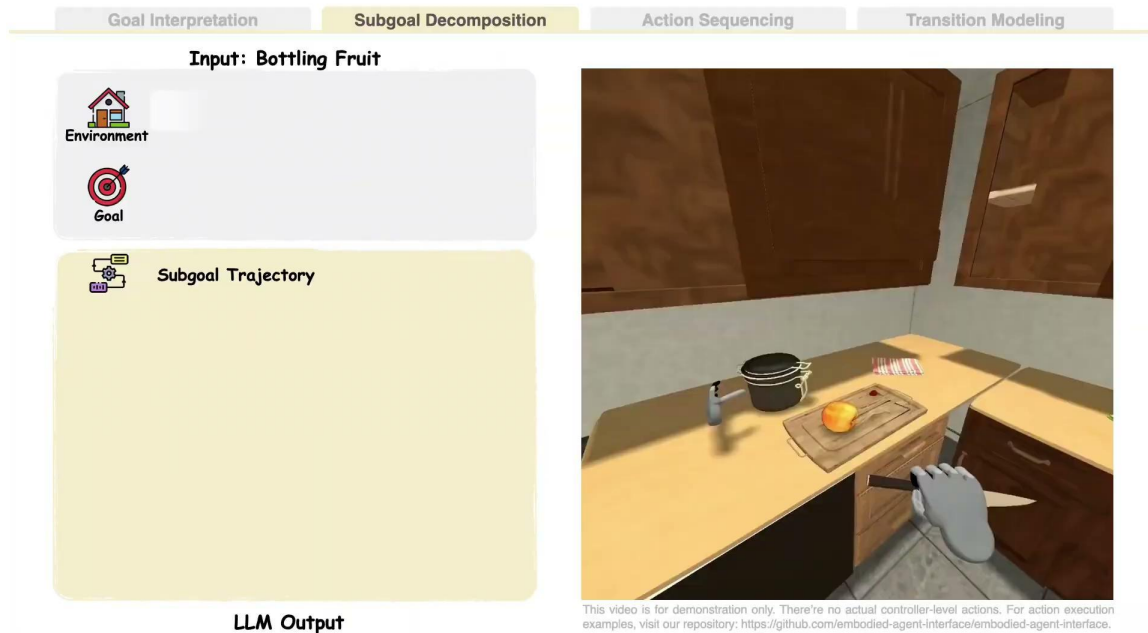


This video is for demonstration only. There's no actual controller-level actions. For action execution examples, visit our repository: https://github.com/embodied-agent-interface/embodied-agent-interface.

# M4. Transition Modeling

Transition Modeling serves as the low-level controller to guide the simulator in performing state transitions from preconditions to post-effects.
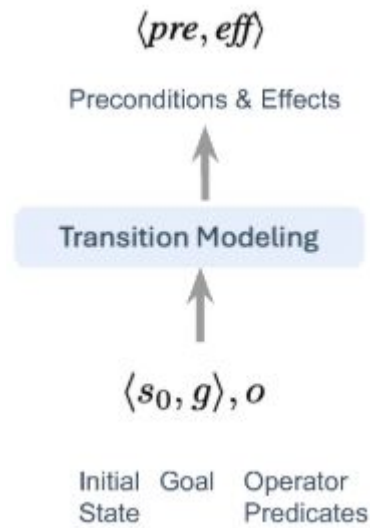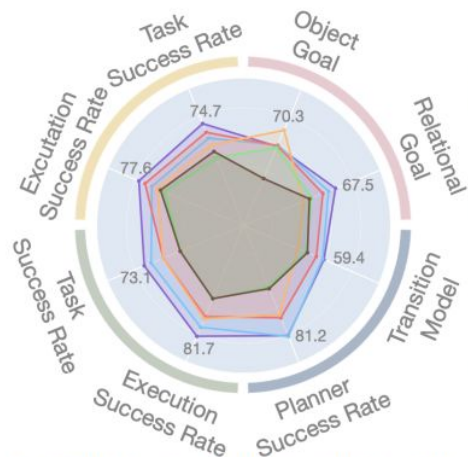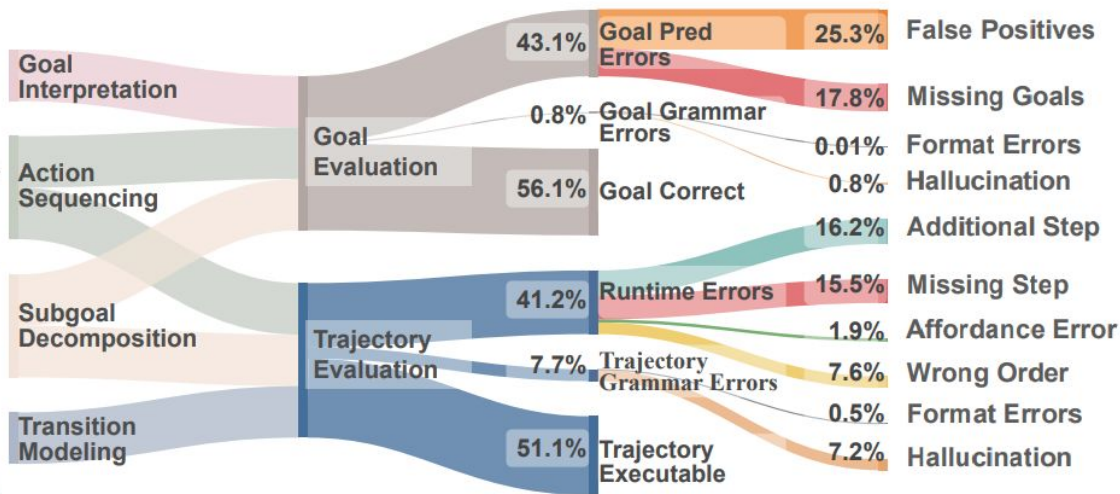
# Evaluation

## Grammar Error

### Parsing

PLACE_ONFLOOR(floor.0)

❌ Unknown action PLACE_ONFLOOR

### Action-Arg Len

GRASP(rag.0, bowl.1)

❌ GRASP only has one param

### Hallucination

RINSE(hand.65)

❌ hand.65 is not in the scene

## Goal Satisfaction Error

### Missing State

**Goal**

*on*(television.410) and
*facing*(agent.65, television.410)

**LLM Output**

. . .
FIND(television.410)
SWITCH_ON(television.410)

**Error Info: State Unsatisfied**

❌ Missing Final State
*facing*(agent.65, television)

### Missing Relation

**Goal**

*next_to*(plywood.78, plywood.79)  and
*next_to*(plywood.79, plywood.80)

**LLM Output**

...LEFT_PLACE_NEXTTO(plywood.79)
LEFT_GRASP(plywood.79)
LEFT_PLACE_NEXTTO(plywood.80)

**Error Info: Relation Unsatisfied**

❌ Missing Final Relation
*next_to*(plywood.78,plywood.79)

### Missing Goal Action

**Goal**

TOUCH(cat)

**LLM Output**

. . .
FIND(cat.1000)
TURN_TO(cat.1000)

**Error Info: Action Unsatisfied**

❌ Missing Goal Action
TOUCH(cat.1000)

## Trajectory – Runtime Error

### Wrong Order

WALK(table.355)
SIT(chair.356)
FIND(novel.1000)
GRAB(novel.1000)

VirtualHome

❌ Precondition
not *sitting*(agent.65) = False

✔ Historical State
not *sitting*(agent.65) = True

### Missing Step

. . .
CLOSE(fridge.0)
SLICE(strawberry.0)
SLICE(peach.0)

BEHAVIOR

❌ Precondition
*holding*(knife.0) = *False*

❌ Historical State
*holding*(knife.0) = *False*

### Affordance Error

LEFT_RELEASE
OPEN(shelf.16)

LEFT_RELEASE
LEFT_GRASP(pool.50)

BEHAVIOR

❌ Precondition
shelf.16 *not openable*

❌ Precondition
pool.50 *not grabbable*

### Additional Step

OPEN(top_cabinet.27)
RIGHT_GRASP(soap.79)
. . .
OPEN(top_cabinet.27)

BEHAVIOR

❌ Current State
*open*(top_cabinet.27) = *True*

⏳ Expected State
*open*(top_cabinet.27) = *False*

# AgentBoard: An Analytical Evaluation Board of Multi-turn LLM Agents

- **Task diversity** is necessary to cover various agent tasks such as embodied, web, and tool agents.
- **Multi-round** interaction is critical to mimic realistic scenarios.
- Evaluating agents in **partially-observable** environments, where they must actively explore to understand their surroundings, is essential for practical assessments.

# AlfWorld

Interactive text environment that require agents to explore surroundings and perform commonsense tasks like "put two soap bars in garbagecan".



Goal: heat some apple and put it in fridge.

**> check valid actions**
Choose an action from these valid actions: go to cabinet 1, go to cabinet 2, go to cabinet 3, …

**> go to fridge 1**
The fridge 1 is closed.

**> open fridge 1**
You open the fridge 1. The fridge 1 is open. In it, you see a apple 1, a bowl 3, a cup 2, a cup 1, a egg 3, a lettuce 1, a potato 2, a potato 1, and a tomato 1. (reward: 0.25)

…

# ScienceWorld

Interactive text environment testing scientific commonsense,
e.g."measure the melting point of the orange juice".



Goal: Your task is to find the animal with the longest life span.

…

**> go to outside**
You move to the outside. (reward: 0.3333333333333333)

**> look around**
This outside location is called the outside. Here you see: the agent a substance called air an axe a baby brown bear ...

**> focus on baby brown bear**
You focus on the baby brown bear. (reward: 0.6666666666666666)

…

# BabyAI

Interactive 20x20 grid environment where agents navigate and interact with objects within a limited sight range.



(a) GoToObj: "go to the blue ball"

(b) PutNextLocal: "put the blue key next to the green ball"

(c) BossLevel: "pick up the grey box behind you, then go to the grey key and open a door". Note that the green door near the bottom left needs to be unlocked with a green key, but this is not explicitly stated in the instruction.

Goal: Open the red door, and open the blue door

**> check valid actions**
You can take the following actions: turn left, turn right, move forward, toggle and go through blue closed door 1, go to blue closed door 1, check available actions
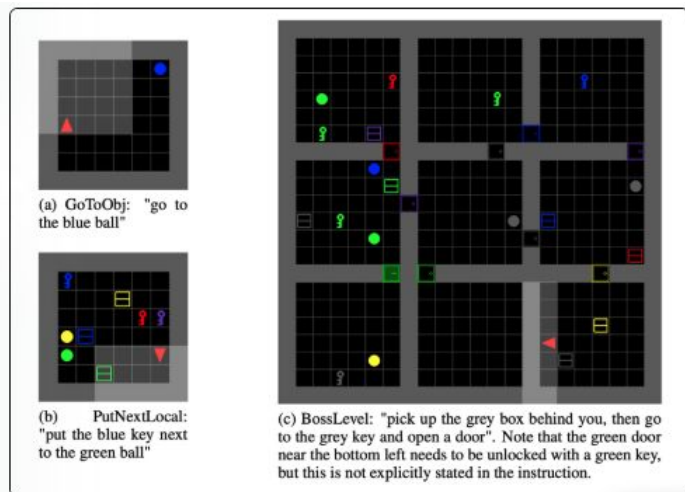
**> go to blue closed door 1**
In front of you in this room, you can see several objects: There is a blue closed door 1 right in front of you 1 steps away. The room has walls around you. You are not carrying anything. (reward: 0.25)
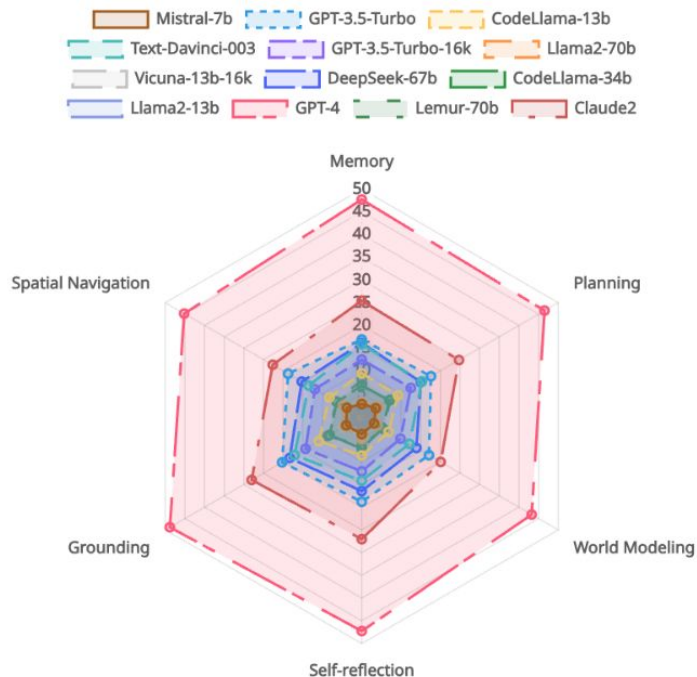
**> toggle and go through blue door 1**
The action is not recognized. Please check valid actions.

…

# Evaluation



| | AlfWorld | ScienceWorld | BabyAI |
|---|---|---|---|
| **Memory** | | | |
| 1. Could finish tasks within 2k tokens | 1 | 2 | 1 |
| 2. Could finish task within 4k tokens | | | |
| 3. Otherwise | | | |
| **Planning** | | | |
| 1. ≤ 3 subgoals on average | 1 | 2 | 2 |
| 2. ≤ 5 subgoals on average | | | |
| 3. Otherwise | | | |
| **World Modeling** | | | |
| 1. Requires no additional knowledge other than instruction | 3 | 3 | 2 |
| 2. Requires knowledge of the environment from exploration | | | |
| 3. Requires commonsense knowledge in addition to knowledge from environment | | | |
| **Self-Reflection** | | | |
| 1. Detailed feedback and error message with instruction for the next step. | 3 | 2 | 2 |
| 2. Not very detailed feedback and error message | | | |
| 3. No error message, e.g. "no change in state" | | | |
| **Grounding** | | | |
| 1. No specific action format is required, could recognize similar actions | 2 | 3 | 2 |
| 2. Action format is required | | | |
| 3. Action format hard to follow | | | |
| **Spatial Navigation** | | | |
| 0. No spatial navigation | 1 | 1 | 1 |
| 1. 2D navigation | | | |