# A holistic and critical look at language agents

Yu Su

The Ohio State University

# Outline

- Language Agents: Introduction

- The Power of Language: Memory, Planning, Grounding, Tool Use

- Future Directions

# Outline

- Language Agents: Introduction


- The Power of Language: Memory, Planning, Tool Use


- Future Directions

# The rise, and the divide

**Bill Gates**

Agents are bringing about the **biggest revolution in computing** since we went from typing commands to tapping on icons.

**Andrew Ng**

I think AI agentic workflows will drive **massive AI progress** this year.

**Sam Altman**

2025 is when **agents will work**.

Current agents are just thin wrappers around LLMs.

Autoregressive LLMs can never reason or plan.

Auto-GPT's limitations in … reveal that it is far from being a practical solution.

# Why agents again?



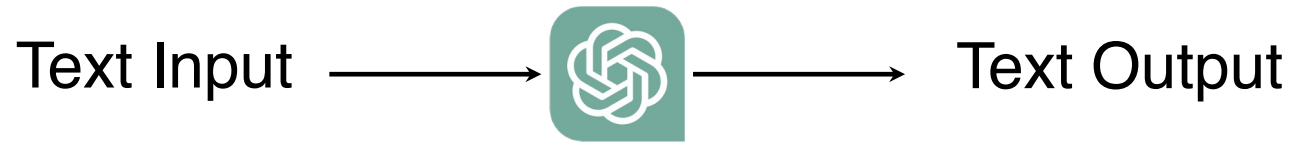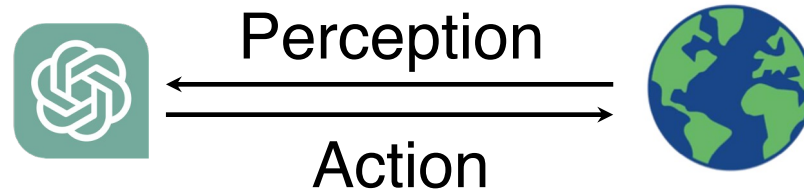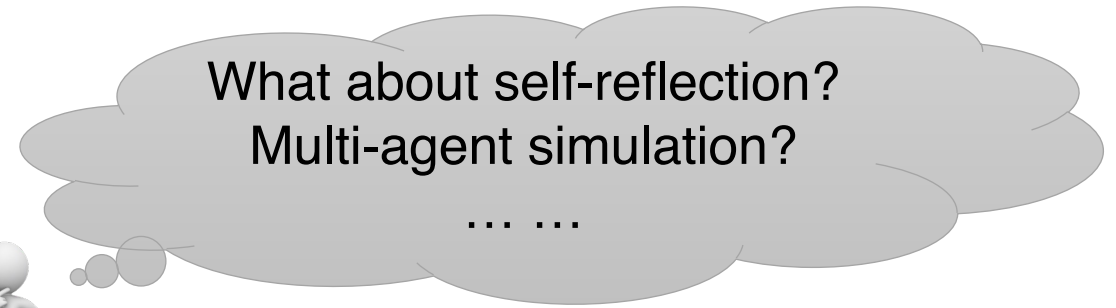"An **agent** is anything that can be viewed as perceiving its **environment** through **sensors** and acting upon that environment through **actuators**"

— Russel & Norvig, *AI: A Modern Approach*

# 'Modern' agent = LLM + external environment?

Text Input → [GPT] → Text Output

**Language Models**

What about self-reflection?
Multi-agent simulation?

… …

Perception
Action

**LLM-based Agents**

# Two competing views

**LLM-first view:** We make an LLM into an agent

- Implications: scaffold on top of LLMs, prompting-focused, heavy on engineering

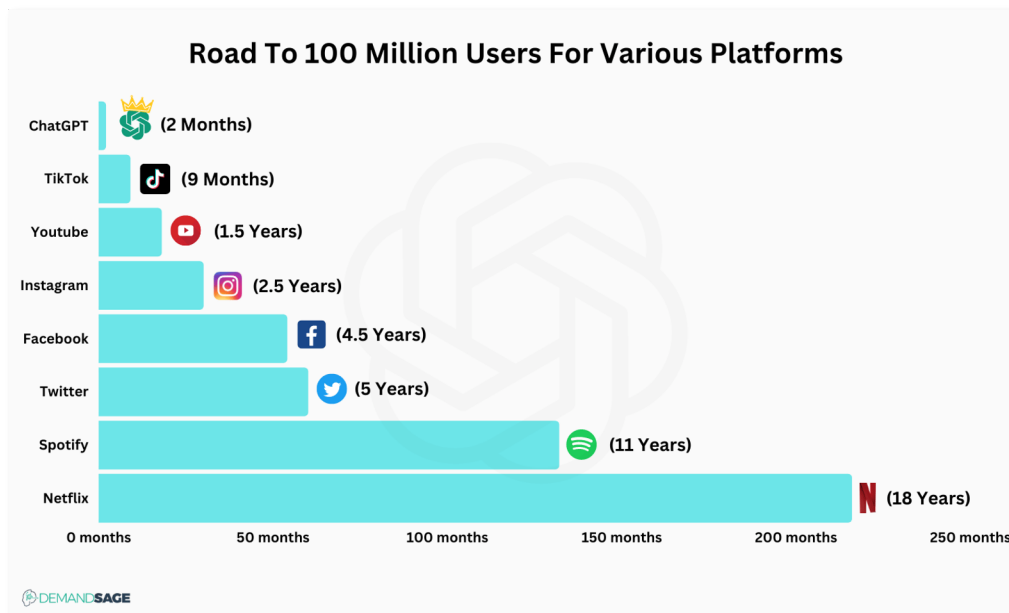**Agent-first view**: An AI agent integrates LLMs so it can use language for reasoning and communication

- Implications: All the same challenges faced by previous AI agents (e.g., *perception, reasoning, world models, planning*) still remain, but we need to re-examine them through the new lens of LLMs and tackle the new challenges and opportunities (e.g., *synthetic data, self-reflection, internalized search*)
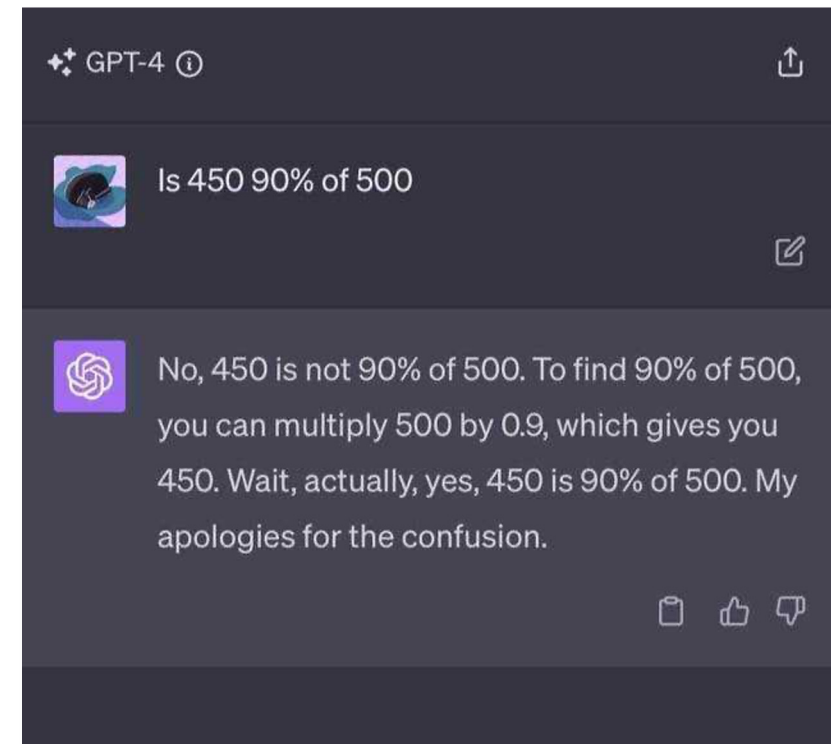
# What's fundamentally different now?

Contemporary AI agents, with integrated LLM(s), can *use language as a vehicle for reasoning and communication*

⬆ Instruction following, in-context learning, output customization

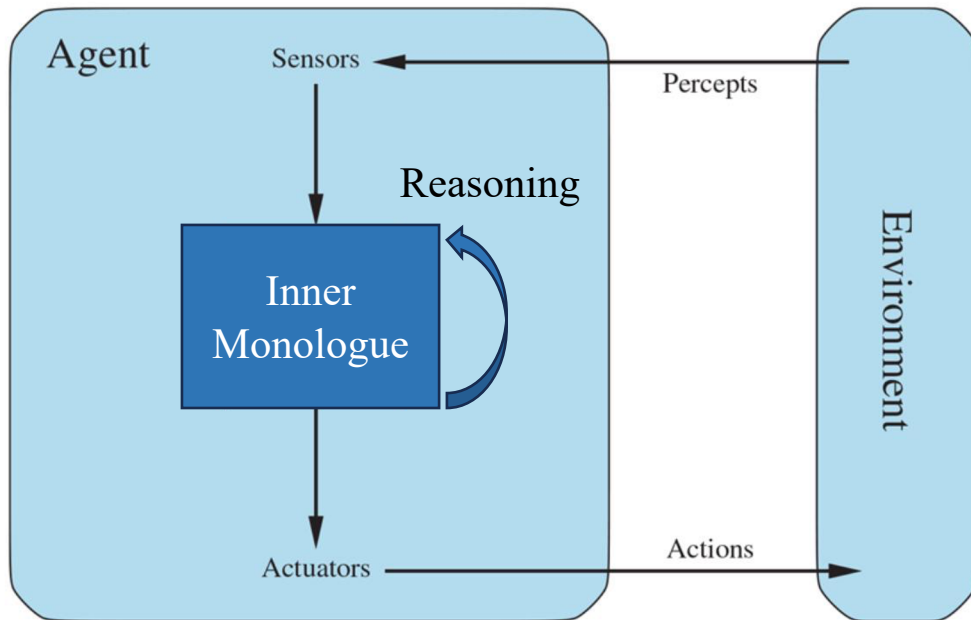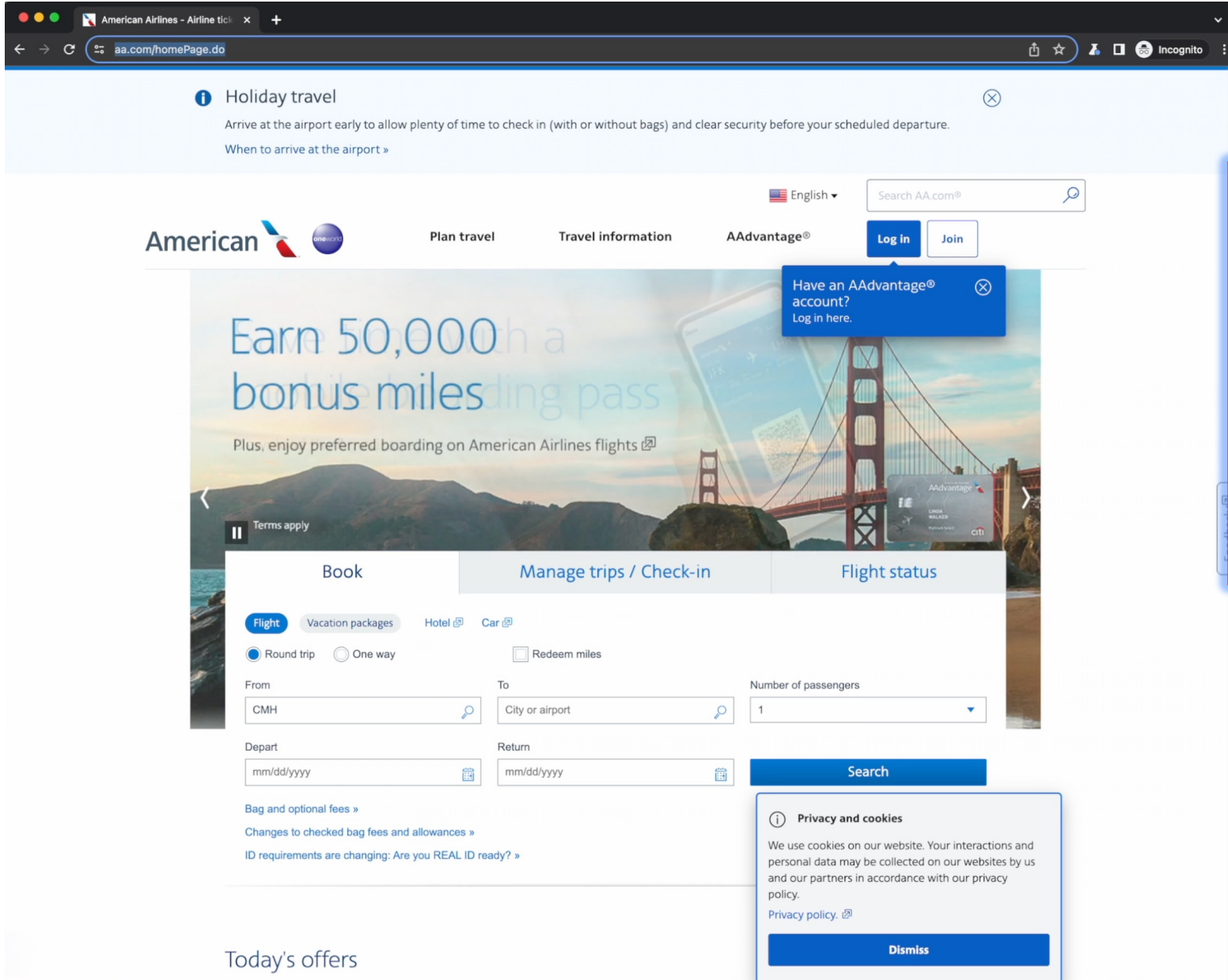⬆ Reasoning (for better acting): state inferences, self-reflection, replanning, etc.

## Road To 100 Million Users For Various Platforms

- ChatGPT (2 Months)
- TikTok (9 Months)
- Youtube (1.5 Years)
- Instagram (2.5 Years)
- Facebook (4.5 Years)
- Twitter (5 Years)
- Spotify (11 Years)
- Netflix (18 Years)

0 months — 50 months — 100 months — 150 months — 200 months — 250 months

DEMANDSAGE

https://www.demandsage.com/chatgpt-statistics/

GPT-4

Is 450 90% of 500

No, 450 is not 90% of 500. To find 90% of 500, you can multiply 500 by 0.9, which gives you 450. Wait, actually, yes, 450 is 90% of 500. My apologies for the confusion.

https://www.reddit.com/r/ChatGPT/comments/16jvl4x/wait_actually_yes/

9

# Schematic illustration: updated



- Reasoning by generating tokens is **a new type of action** (*vs.* actions in external environments)

- **Internal environment**, where reasoning takes place in an inner monologue fashion

- **Self-reflection** is a 'meta' reasoning action (i.e., reasoning over the reasoning process), akin to metacognitive functions

- **Reasoning is for better acting**, by inferring environmental states, retrospection, etc.

- **Percept** and **external action spaces** are substantially expanded, thanks to using language for communication and multimodal perception

# Web agent example



Task: Search for the flight status of AA 3942 leaving on Dec 29

# Let's call them language agents

- These contemporary AI agents capable of using language for reasoning and communication are best called "**language agents**," for language being their most salient trait.

- What about *multimodal agents*?
  - While there's perception of other modalities, language is still doing the heavy lifting (reasoning and communication)

- What about simply *LLM agents*?
  - The key is using language for reasoning and communication, but that doesn't have to come from an LLM; that may turn out to be a means to an end
  - Maybe in a few years, we will move beyond LLMs, but the need for universal language understanding and production in agents will remain

# Evolution of AI agents



**Logical Agent**

**Neural Agent**

**Language Agent**

| | | | |
|---|---|---|---|
| **Expressiveness** | | | |
| **Reasoning** | | | |
| **Adaptivity** | | | |

# Evolution of AI agents



**Logical Agent**  **Neural Agent**  **Language Agent**

| | Logical Agent | Neural Agent | Language Agent |
|---|---|---|---|
| **Expressiveness** | Low<br>bounded by the logical language | | |
| **Reasoning** | Logical inferences<br>sound, explicit, rigid | | |
| **Adaptivity** | Low<br>bounded by knowledge curation | | |

# Evolution of AI agents



**Logical Agent**



**Neural Agent**



**Language Agent**

| | Logical Agent | Neural Agent | Language Agent |
|---|---|---|---|
| **Expressiveness** | Low<br>bounded by the logical language | Medium<br>anything a (small-ish) NN can encode | |
| **Reasoning** | Logical inferences<br>sound, explicit, rigid | Parametric inferences<br>stochastic, implicit, rigid | |
| **Adaptivity** | Low<br>bounded by knowledge curation | Medium<br>data-driven but sample inefficient | |

# Evolution of AI agents



**Logical Agent**    **Neural Agent**    **Language Agent**

| | Logical Agent | Neural Agent | Language Agent |
|---|---|---|---|
| **Expressiveness** | Low<br>bounded by the logical language | Medium<br>anything a (small-ish) NN can encode | High<br>almost anything, esp. those verbalizable |
| **Reasoning** | Logical inferences<br>sound, explicit, rigid | Parametric inferences<br>stochastic, implicit, rigid | Language-based inferences<br>fuzzy, semi-explicit, flexible |
| **Adaptivity** | Low<br>bounded by knowledge curation | Medium<br>data-driven but sample inefficient | High<br>strong prior from LLMs + language use |

# A conceptual framework for language agents

# Outline

- Language Agents: Introduction

- The Power of Language: Memory, Planning, Grounding, Tool Use

- Future Directions

# Memory

# Long-term memory in LLMs

• LLMs still rely on RAG.

**ChatGPT**

**You**

Who won the 2024 Super Bowl?

**ChatGPT**

The Super Bowl in 2024 was held on February 11, 2024. To find the most current information about the winner, I can browse the internet for you. Let me do that now.

The Kansas City Chiefs won the 2024 Super Bowl against the San Francisco 49ers with a final score of 25-22. The game went into overtime, and the Chiefs secured the win with a touchdown just seconds before the end of the game (Yahoo Sports) (KCUR) (PinkNews) (Just Jared).

# Current RAG does not always work

Passages

Which 🏛 Stanford professor works on the
neuroscience of 🎗 Alzheimer's?

Thomas

Answer:

# Current RAG does not always work

# Current RAG does not always work

# Long-term memory in humans

- One well-established theory of human long-term memory is the **hippocampal indexing theory** (Teyler et al. 1986).
  - The **hippocampus (blue)** is a store for **indices** (which point to memories stored in **neocortex**) and **associations** between them.
  - Memories can be anything (i.e., people, events, places, etc.)

# Long-term memory in humans

- Indexing procedure enables two fundamental faculties of human memory:

  - **Pattern separation:** process for differentiating memories (neocortex and parahippocampus)
  - **Pattern completion:** process for recovering complete memories from relevant associations (mostly hippocampus, specifically CA3)

# HippoRAG:
# Neurobiologically-inspired LTM for LLMs

# HippoRAG:
# Neurobiologically-inspired LTM for LLMs

# HippoRAG & the three components of LTM

**Neocortex**
Perception, linguistic
abilities & reasoning

**Parahippocampus**
Bridge between areas;
working memory

**Hippocampus**
Indexing & auto-
associative memory

# HippoRAG & the three components of LTM



|  | | |
|---|---|---|
| Neocortex | Parahippocampal Regions | Hippocampus |
| LLM | Retrieval Encoders | KG + Personalized PageRank |

**Offline Indexing**

Passages → Open IE:
(Thomas, *researches*, Alzheimer's)
(Stanford, *employs*, Thomas)

# HippoRAG & the three components of LTM



| Neocortex | Parahippocampal Regions | Hippocampus |
|---|---|---|
| LLM | Retrieval Encoders | KG + Personalized PageRank |

**Offline Indexing**

Passages → Open IE → (Thomas, *researches*, Alzheimer's) (Stanford, *employs*, Thomas)

**Online Retrieval**

Query → NER → Stanford / Alzheimer's

Node Specificity

# HippoRAG is a SoTA memory retriever



Average Multi-Hop QA Performance

32

# Practical aspects of HippoRAG

- Better retrieval directly translates to better QA (avg. 5% better)

- With *single-step multi-hop retrieval*, it's 10-30 times cheaper and 6-13 times faster than iterative RAG like IRCoT with similar or better recall

- Worried about indexing cost? Llama-3-8B achieves similar performance as GPT-3.5

# Planning



Core Competencies

Multi-agent/Theory of Mind — Continual Learning

Grounding — Planning — Tool Use

Reasoning — World Models

Multimodal Perception — Memory — Embodiment

# *Is Your LLM Secretly a World Model of the Internet?* Model-based Planning for Web Agents

**Yu Gu**\*, **Boyuan Zheng**\*, Boyu Gou, Kai Zhang,

Cheng Chang, Sanjari Srivastava, Yanan Xie, Peng Qi,

Huan Sun, Yu Su

NLP    rby AI

# Planning paradigms for language agents



(a) reactive
- 😊 fast, easy to implement
- 😐 greedy, short-sighted

(b) tree search with real interactions
- 😊 systematic exploration
- 😐 irreversible actions, unsafe, slow

(c) model-based planning
- 😊 faster, safer, systematic exploration
- 😐 how to get a world model?

36

# What's ... a world model?

A computational model of environment transition dynamics

$$\hat{T} : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$$

*If I do this ($a_t$) right now ($s_t$), what would happen next ($s_{t+1}$)?*

# Why hasn't it been done already?



VS



And **billions of other websites** on the Internet!

# LLMs can predict state transitions



The page will navigate to a detailed product page for the "Mens Flowers Casual Aloha Hawaiian Shirt Summer Short Sleeve Beach T-Shirt Regular Fit Button Down Dress Shirts." This new page will likely contain additional information about the product including more detailed specifications, customer reviews, larger images, sizing options, and possibly a larger "Add to Cart" button. Other elements from the current category view like the grid of products will be replaced with the detailed view of this specific product.

# WebDreamer: model-based planner for web agents



Please navigate to the 'Data Storage' category and purchase the least expensive disk with 512GB of storage.

# WebDreamer: model-based planner for web agents

Please navigate to the 'Data Storage' category and purchase the least expensive disk with 512GB of storage.

# WebDreamer: model-based planner for web agents

Please navigate to the 'Data Storage' category and purchase the least expensive disk with 512GB of storage.



③ Type 'Disk'

① Click 'Office Products'

② Click 'Electronics'

The 'Office Products' category will display three sub-categories: 'Office Electronics', 'Office & School Supplies', and 'Office Furniture & Lighting'.

Click 'Office Electronics'

The webpage will display 'Office Electronics' sub-category results with products, and the sub-menu will show 'Printers&Accessories' and other categories.

# WebDreamer: model-based planner for web agents

**Please navigate to the 'Data Storage' category and purchase the least expensive disk with 512GB of storage.**



③ Type 'Disk'

① Click 'Office Products'

② Click 'Electronics'

**Stage I: Simulation**

① The 'Office Products' category will display three sub-categories: 'Office Electronics', 'Office & School Supplies', and 'Office Furniture & Lighting'.

Click 'Office Electronics'

The webpage will display 'Office Electronics' sub-category results with products, and the sub-menu will show 'Printers&Accessories' and other categories.

v = 0.4

# WebDreamer: model-based planner for web agents

Please navigate to the 'Data Storage' category and purchase the least expensive disk with 512GB of storage.



## Stage I: Simulation

**1** The 'Office Products' category will display three sub-categories: 'Office Electronics', 'Office & School Supplies', and 'Office Furniture & Lighting'.

*Click 'Office Electronics'* → The webpage will display 'Office Electronics' sub-category results with products, and the sub-menu will show 'Printers&Accessories' and other categories. **v = 0.4**

**2** The 'Electronics' category will display three sub-categories: 'Computers & Accessories', 'Accessories & Supplies', and 'Car & Vehicle Electronics'.

*Click 'Computer & Accessories'* → The webpage will display 'Computer Accessories' sub-category results, including 'Data Storage', 'Tablet Accessories', and others. **v = 0.8**

**3** The webpage will display search results, including a list of products, each of which includes the product title, price, and an 'Add to Cart' button.

*Click 'Electronics'* → The 'Electronics' category will display three sub-categories: 'Computers & Accessories', 'Accessories & Supplies', and 'Car & Vehicle Electronics'. **v = 0.1**

## Stage II: Execution



44

# Results on VisualWebArena

Model-based planning is more accurate than reactive planning and more efficient than tree search



Baselines are from Koh et al., (2024)

# Grounding

# *Navigating the Digital World as Humans Do*: Universal Visual Grounding for GUI Agents

**Boyu Gou**, Ruohan Wang, Boyuan Zheng, Yanan Xie,

Cheng Chang, Yiheng Shu, Huan Sun, Yu Su

# Embodiment of Web/GUI Agents: Evolution



Xiang Deng

| | Mind2Web (23' May) | | |
|---|---|---|---|
| **Sensory Inputs** | HTML/DOM | | |
| **Effectors** | Multi-choice Selection | | |

# Embodiment of Web/GUI Agents: Evolution



## Action Description

Move the cursor over the **"Find Your Truck" button** located in the central portion of the webpage, just below the input fields for rental details, and perform **a click action.**

A: <a id="0">Moving Trucks & Accessories</a>
B: <input type="text" id="1">placeholder="US City,State or Zip Code"</input>
... ... ...
F: <input type="radio" id="5">No name="one-way-radio"</input>
G: <input type="button" id="6">value="Find Your Truck"</input>
H: None

### Element Attributes
TEXT: Find Your Truck
TYPE: BUTTON

Boyuan Zheng

### Image Annotation
CHOICE: G

Boyu Gou

### Textual Choices
CHOICE: G

| | Mind2Web (23' May) | SeeAct (24' Jan) |
|---|---|---|
| **Sensory Inputs** | HTML/DOM | Screenshot + DOM |
| **Effectors** | Multi-choice Selection | Multi-choice Selection |

# Embodiment of Web/GUI Agents: Evolution



**Text-based Representations: Limitations**

- *Noisy and incomplete*
  - **95.9%** of home pages have accessibility conformance errors
  - Avg. **56.8** errors per page[1]
- *Additional input increases latency and inference costs*
  - Consuming more tokens
  - Difficult and time-consuming to get
  - Compounding over long horizon

| | Mind2Web (23' May) | SeeAct (24' Jan) | |
|---|---|---|---|
| **Sensory Inputs** | HTML/DOM | Screenshot + DOM | |
| **Effectors** | Multi-choice Selection | Multi-choice Selection | |

[1] https://webaim.org/projects/million/

# Embodiment of Web/GUI Agents: Evolution

**SeeAct-V: Human-like, Vision-centric GUI Agent**

Boyu Gou

**Vision-Only Observation**

TASK: Find the cheapest 4k monitor

**Planning**

Element Description:
**The search bar at the top of the page**
**Action:** Type
**Value:** 4k monitor

**Grounding**

What are the pixel coordinates of the element corresponding to "..."?

**(556, 26)**

**Human-like Operation**

**Click**(556, 26)
**Type**("4k monitor")

| | Mind2Web (23' May) | SeeAct (24' Jan) | SeeAct-V (24' Aug) |
|---|---|---|---|
| **Sensory Inputs** | HTML/DOM | Screenshot + DOM | Screenshot Only |
| **Effectors** | Multi-choice Selection | Multi-choice Selection | Pixel-level Operations |

Also see CogAgent and SeeClick for prior work along this line

# Visual Grounding is the Bottleneck

## GPT-4V(ision) is a Generalist Web Agent, if Grounded

**Boyuan Zheng** [1]  **Boyu Gou** [1]  **Jihyung Kil** [1]  **Huan Sun** [1]  **Yu Su** [1]

https://osu-nlp-group.github.io/SeeAct

## Abstract

The recent development on large multimodal models (LMMs), especially GPT-4V(ision) and Gemini, has been quickly expanding the capability boundaries of multimodal models beyond traditional tasks like image captioning and visual question answering. In this work, we explore the potential of LMMs like GPT-4V as a generalist web agent that can follow natural language instructions to complete tasks on any given website. We propose SEEACT, a generalist web agent that harnesses the power of LMMs for integrated visual understanding and acting on the web. We evaluate on the recent MIND2WEB benchmark. In addition to standard offline evaluation on cached websites, we enable a new online evaluation setting by developing a tool that allows running web agents on live websites. We show that GPT-4V presents a great potential for web agents—it can successfully complete 51.1% of the tasks on live websites if we manually ground its textual plans into

Figure 1: SEEACT leverages an LMM like GPT-4V to visually perceive websites and generate plans in textual forms. The textual plans are then grounded onto the HTML elements and operations to act on the website.

# Referring Expressions for GUIs are Diverse



1. Red icon labeled "UNIQLO"
2. Button at the top left corner
3. Navigate back to the homepage

1. Hollow heart button
2. Button below the Pokémon shirt
3. Favor the Pokémon shirt

- *Visual Referring Expressions*
  - Salient visual features like textual content, element type (button, input field, checkbox, etc.), shape, color, …

- *Positional Referring Expressions*
  - including **absolute** (e.g., "*at the top left of the page*") and **relative** positions (e.g., "*to the right of element X*")

- *Functional Referring Expressions*
  - Referring to elements by their functions

- *Hybrid*
  - *"click the heart button under the Pokemon shirt to add it to favorite."*

# Shared Designs across GUIs



*"Go to homepage"*

*"Go to homepage"*

*"Open Maps"*

*"iPhone 16"*

# UGround: Universal Visual Grounding for GUI Agents

- What do we mean by universal?
  - Generalize across different web, desktop, and mobile environments
  - Plug and play with different planners / LLMs
  - Handle input images of different resolutions
- Modeling
  - **Input**: Screenshot + RE; **Output**: Pixel Coordinates
  - LLaVA backbone: CLIP@224px, Vicuna-1.5-7B-16K, AnyRes

# Synthetic Data is Key

- Collect all relevant information about an element
  - Visual (i.e., screenshot)
  - HTML attributes like `inner-text`, `alt-text` as additional visual cues
  - Accessibility labels like `alt`, `aria-label` as functionality cues

- Generate referring expressions using (multimodal) LLMs and rules

- + Positional information

**Screenshot**



**HTML**

**type**: button
**aria-label**: menu

🌋 **MLLM Interpretation**

*The image shows a familiar "hamburger menu" icon, which consists of three horizontal, evenly spaced lines stacked vertically. This icon is often used in web and mobile interfaces to indicate a collapsible or expandable menu, commonly referred to as the "menu" or "navigation" icon. Its minimalistic design helps users quickly recognize it as a toggle for hidden menu items or additional navigation options. Given the context from the HTML information, it appears this icon is intended to open a menu when clicked.*

🦙 **LLM Rewrite**

*three-line menu icon*

📐 **Rules**

**Absolute Position**: Top Left Corner
**Relative Position**: to the left of "Donate"

…

🧩 **Final Referring Expression**

*three-line menu icon, at the top left corner of the page*

# Training Data Statistics

| Dataset | Annotation | # of Elements | # of Screenshots | Platform |
|---|---|---|---|---|
| Web-Hybrid (Ours) | Rule + LLM | 9M | 773K | Web |
| Web-Direct (Ours) | GPT | 408K | 408K | Web |
| GUIAct (Chen et al., 2024) | GPT + Human | 140K | 13K | Web |
| AndroidControl (Li et al., 2024b) | Human | 47K | 47K | Android |
| Widget Caption (Li et al., 2020b) | Human | 41K | 15K | Android |
| UIBert (Bai et al., 2021) | Human | 16K | 5K | Android |
| AITZ (Zhang et al., 2024b) | GPT + Human | 8K | 8K | Android |
| Total | | 10M | 1.3M | Web + Android |

# Human-like Vision-centric Embodiment is the Future for Web/GUI Agents

- Most comprehensive evaluation of GUI Agents to date
- SeeAct-V + UGround outperforms prior SoTA despite its minimalist design



Figure 1: Examples of agent tasks across platforms and performance on **GUI grounding** (♣: ScreenSpot), **offline agent** (♠: Multimodal-Mind2Web, AndroidControl, and OmniAct), and **online agent benchmarks** (♥: Mind2Web-Live and AndroidWorld) when using GPT-4 as the planner.

# Error Analysis

- Most remaining errors are due to planning rather than grounding
- Grounding errors on mobile and desktop are higher, but that doesn't necessarily manifest in the agent performance

# Tool Use

# LLMs in the Imaginarium:
# Tool Learning through Simulated Trial and Error

**Boshi Wang**, Hao Fang, Jason Eisner,

Benjamin Van Durme, Yu Su

# We don't master a tool by just reading the 'manual'



this too and humans once they reach about age five

https://www.youtube.com/watch?v=s2IBayVsbz8

- Trial and error

- Imagination

- Memory

# Tool learning through *simulated trial and error*

### Synthesizing diverse **perception-decision-execution** data with LLMs

# STE substantially improves tool use

**ICL**: Retrieve from STE synthesized examples. **Fine-Tuning**: Fine-tune LLM with STE examples



65

# A few remarks

- Language agents are still AI agents. The same challenges remain. But the new capability of using language for reasoning and communication allows us to re-examine these challenges through a new, powerful lens.

- Natural selection is a real magical power. Over billions of years of evolution, natural selection has found numerous good designs ('local optima') that *just make sense*; we should take some inspirations!

- Tons of new venues for interesting research. We are barely scratching the surface…

# Outline

- Language Agents: Introduction

- The Power of Language: Memory, Planning, Grounding, Tool Use

- Future Directions

# We are just standing at the dawn of a long journey

# Future directions for language agents



- *Hierarchical planning?*

- *How far can (tree) search go?*

- *How to evaluate (partial) plans? Value functions? Reward models?*

- *Speculative planning?*

Chen et al., *When is Tree Search Useful for LLM Planning? It Depends on the Discriminator*, ACL 2024
Gu et al., *Is Your LLM Secretly a World Model of the Internet? Model-based Planning for Web Agents*, arXiv 2024

# Future directions for language agents



Cross-cutting Issues
- Safety
- Evaluation
- Synthetic Data
- Efficiency
- Applications

Core Competencies
- Multi-agent / Theory of Mind
- Continual Learning
- Grounding
- Planning
- Tool Use
- Reasoning
- World Models
- Multimodal Perception
- Memory
- Embodiment

Environments
- Humans
- Agents
- Databases
- Web
- Physical World
- Computers
- APIs
- Documents

- Agents need to learn **perception-decision-execution** capabilities

- Data on the Internet is mostly artifacts from such processes, not capturing the processes *per se*

- Synthesizing data with LLMs provides a possibility to uncover (some of) these hidden processes

Wang et al., *LLMs in the Imaginarium: Tool Learning through Simulated Trial and Error*, ACL 2024
Gou et al., *Navigating the Digital World as Humans Do: Universal Visual Grounding for GUI Agents*, arXiv 2024

# Future directions for language agents



- Existing MLLMs are not optimized for GUI understanding.
  - Resolution/grounding/fine-grained understanding

- Current MLLMs do not seem to fuse visual and language well. New arch or training recipes are needed.

# Future directions for language agents



Cross-cutting Issues

| Safety | Evaluation | Synthetic Data | Efficiency | Applications |

Core Competencies

Multi-agent / Theory of Mind · Continual Learning

Grounding · Planning · Tool Use

Reasoning · World Models
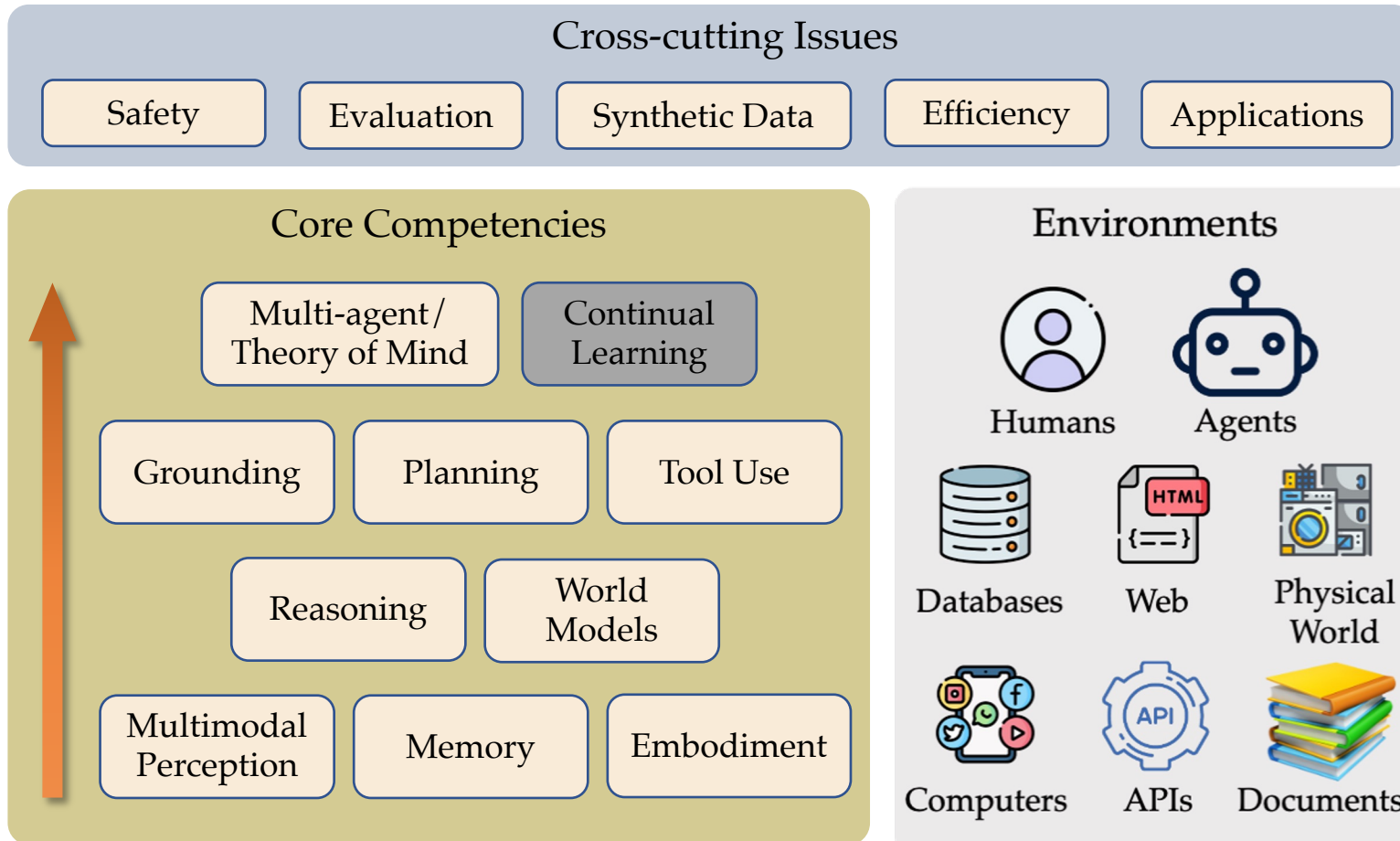
Multimodal Perception · Memory · Embodiment

Environments

Humans · Agents

Databases · Web · Physical World

Computers · APIs · Documents

- Currently, the field is transitioning from *prompting* to *behavior cloning / supervised fine-tuning*

- Behavior cloning is unlikely to be sufficient for generalist agents; they need to explore the environments and *trial and error*

- Challenges from open action space, reward model, and safety

# Future directions for language agents

## Cross-cutting Issues

Safety | Evaluation | Synthetic Data | Efficiency | Applications

## Core Competencies

Multi-agent/ Theory of Mind | Continual Learning

Grounding | Planning | Tool Use

Reasoning | World Models

Multimodal Perception | Memory | Embodiment

## Environments

Humans | Agents

Databases | Web | Physical World

Computers | APIs | Documents

- Agent safety research is far behind agent development and deployment

- Language agents inherent all the safety risks of LLMs (e.g., *bias, fairness, hallucination, privacy, transparency*), amplify some of them (e.g., *workforce displacement*), and bring many new ones (e.g., *irreversible actions*)

# Language Agents: Foundations, Prospects, and Risks

Yu Su, Diyi Yang, Shunyu Yao, Tao Yu

https://language-agent-tutorial.github.io/

# Acknowledgements

# Thanks &