# From Multimodal LLM to Human-level AI

## *Modality*, *Instruction*, *Reasoning*, *Efficiency* and **Beyond**
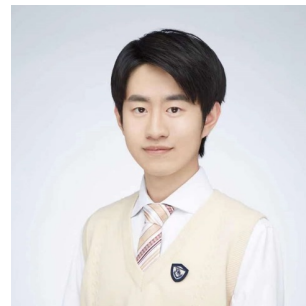
**https://mllm2024.github.io/CVPR2024/**

1

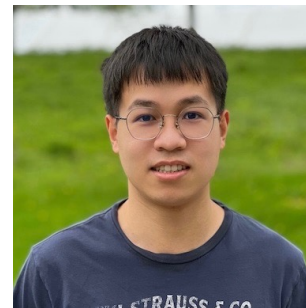**Hao Fei**
*National University of Singapore*

**Yuan Yao**
*National University of Singapore*

**Ao Zhang**
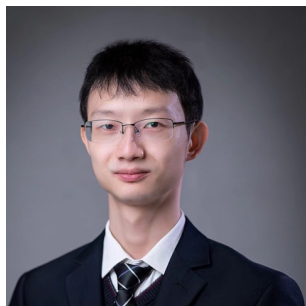*National University of Singapore*

**Haotian Liu**
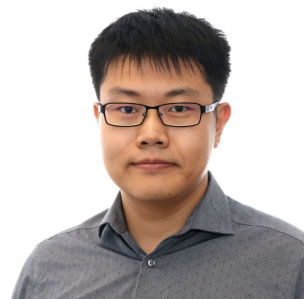*University of Wisconsin-Madison*

**Fuxiao Liu**
*University of Maryland, College Park*

**Zhuosheng Zhang**
*Shanghai Jiao Tong University*

**Hanwang Zhang**
*Nanyang Technological University*

**Shuicheng Yan**
*Kunlun 2050 Research, Skywork AI*

# Part-VI

# Multimodal Reasoning

**Zhuosheng Zhang**
**Tenure-Track Assistant Professor**
*Shanghai Jiao Tong University*

https://bcmi.sjtu.edu.cn/~zhangzs/

# Table of Content

- Basics of Multimodal Reasoning

    × Background, Definition, and Development

- Multimodal Chain-of-Thought Reasoning

    × Paradigm Shift, the Role of Multimodal CoT

- Towards Multimodal LLM Agents

    × Taxonomy, Architecture, Applications

- Challenges

    × Evolutionary Reasoning, Interactive Reasoning, Reasoning Alignment

# 1

# Basics of Multimodal Reasoning
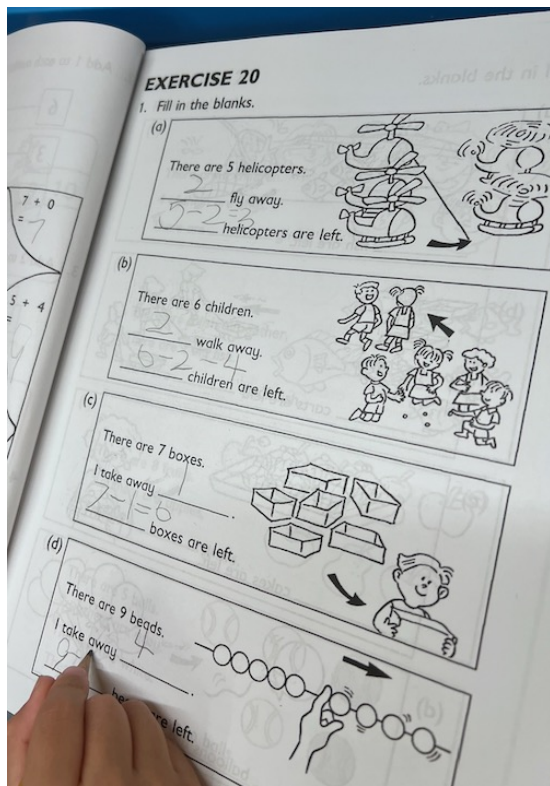
# Multimodal Reasoning



**Prompt**

Please create a blog post about a trip to New York, where a dog and his owner had lots of fun. Include a few pictures of the dog posing happily at different landmarks.

**Model Response**

Hi, I'm Cooper. I'm a 2 year old Goldendoodle who loves to travel. Last month, my family and I went to New York City. It was my first time in the Big Apple, and I had a blast!

We stayed in a hotel near Times Square. It was so much fun to see all the bright lights and people. We went to see a Broadway show, and I even got to meet some of the actors.

One of my favorite parts of New York City was Central Park. It was so big and beautiful. I loved running around and playing fetch. We also went to the Statue of Liberty and the Empire State Building. I was so excited to see all the famous landmarks.

I had so much fun in New York City. It's definitely one of my favorite places to travel. I can't wait to go back someday!

# Multimodal Reasoning

❑ **Concept: a process of deriving high-level conclusions from multiple modalities, possibly via multiple logical steps based on atomic evidences**

---

## Accounting : Investment

**Question:** Consider the three stocks in the following table. Pt represents price at time t, and Qt represents shares outstanding at time t. Stock C splits two for one in the last period. Calculate the rate of return on a price-weighted index of the three stocks for the first period ( t 0 to t 1).

**Option:**
(A) 4.17%
(B) 4.03%
(C) 3.89%

|   | $P_0$ | $Q_0$ | $P_1$ | $Q_1$ | $P_2$ | $Q_2$ |
|---|---|---|---|---|---|---|
| A | 90 | 100 | 95 | 100 | 95 | 100 |
| B | 50 | 200 | 45 | 200 | 45 | 200 |
| C | 100 | 200 | 110 | 200 | 55 | 400 |

**GPT-4V:** To calculate the rate of return on a price-weighted index of the three stocks for the first period (t0 to t1), follow these steps: Calculate the initial index value at t0. Calculate the index value at t1. Compute the rate of return using the two index values.

Initial Index Value at t0:
Price-weighted index value at t0 = $(P0(A) + P0(B) + P0(C)) / 3$
= $(90 + 50 + 100) / 3$
= 80
Index Value at t1:
Price-weighted index value at t1 = $(P1(A) + P1(B) + P1(C)) / 3$
= $(95 + 45 + 110) / 3$
= 83.33
Rate of Return:
Rate of Return = (Index Value at t1 - Index Value at t0) / Index Value at t0
= (83.33 − 80) / 80
= 0.041625
= 4.1625%
Therefore, the rate of return on a price-weighted index of the three stocks for the first period is approximately 4.17%.
The correct answer is:
**(A) 4.17%**

**Ground Truth: (A) 4.17%**
**Explanation:** At t = 0, the value of the index is: (90 + 50 + 100)/3 = 80. At $t$ = 1, the value of the index is: (95 + 45 + 110)/3 = 83.333. The rate of return is: (83.333/80) − 1 = 4.17%.

---

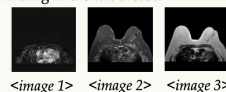### Comprehensive Disciplines

Engineering (26%)  Art & Design (11%)

Business (14%)

Science (23%)   Humanities & Social Sci. (9%)

Medicine (17%)

### Heterogeneous Image Types

Diagrams, Tables, Plots and Charts, Photographs, Chemical Structures, Paintings, Medical Images, Sheet Music, Geometric, Pathology images, Microscopic Images, Comics, …

### Interleaved Text and Images

**Question:** You are shown subtraction <image 1>, T2 weighted <image 2> and T1 weighted axial <image 3> from a screening breast MRI. What is the etiology of the finding in the left breast?

### Expert-level Skills Test

Expert-level Visual Perception

Perception

Knowledge | Reasoning

Domain Expertise, World, Linguistic, Visual Knowledge,…

Logical, Spatial Commonsense, Mathematical,…

---

### Art & Design

**Question:** Among the following harmonic intervals, which one is constructed incorrectly?

**Options:**
(A) Major third <image 1>
(B) Diminished fifth <image 2>
(C) Minor seventh <image 3>
(D) Diminished sixth <image 4>

**Subject:** Music; **Subfield:** Music; **Image Type:** Sheet Music; **Difficulty:** Medium

### Business

**Question:** …The graph shown is compiled from data collected by Gallup <image 1>. Find the probability that the selected Emotional Health Index Score is between 80.5 and 82?

**Options:**
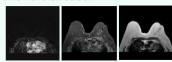(A) 0      (B) 0.2142
(C) 0.3571   (D) 0.5

**Subject:** Marketing; **Subfield:** Market Research; **Image Type:** Plots and Charts; **Difficulty:** Medium

### Science

**Question:** <image 1> The region bounded by the graph as shown above. Choose an integral expression that can be used to find the area of R.

**Options:**
(A) $\int_0^{1.5}[f(x) - g(x)]dx$
(B) $\int_0^{1.5}[g(x) - f(x)]dx$
(C) $\int_0^{2}[f(x) - g(x)]dx$
(D) $\int_0^{2}[g(x) - x(x)]dx$

**Subject:** Math; **Subfield:** Calculus; **Image Type:** Mathematical Notations; **Difficulty:** Easy

---

### Health & Medicine

**Question:** You are shown subtraction <image 1>, T2 weighted <image 2> and T1 weighted axial <image 3> from a screening breast MRI. What is the etiology of the finding in the left breast?

**Options:**
(A) Susceptibility artifact
(B) Hematoma
(C) Fat necrosis   (D) Silicone granuloma

**Subject:** Clinical Medicine; **Subfield:** Clinical Radiology; **Image Type:** Body Scans: MRI, CT.; **Difficulty:** Hard

### Humanities & Social Science

**Question:** In the political cartoon, the United States is seen as fulfilling which of the following roles? <image 1>

**Option:**
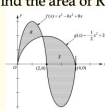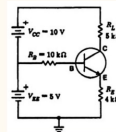(A) Oppressor
(B) Imperialist
(C) Savior   (D) Isolationist

**Subject:** History; **Subfield:** Modern History; **Image Type:** Comics and Cartoons; **Difficulty:** Easy

### Tech & Engineering

**Question:** Find the VCE for the circuit shown in <image 1>. Neglect VBE

**Answer:** 3.75

**Explanation:** …IE = [(VEE) / (RE)] = [(5 V) / (4 k-ohm)] = 1.25 mA; VCE = VCC - IERL = 10 V - (1.25 mA) 5 k-ohm; VCE = 10 V - 6.25 V = 3.75 V

**Subject:** Electronics; **Subfield:** Analog electronics; **Image Type:** Diagrams; **Difficulty:** Hard

Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y. and Wei, C. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. CVPR 2024.

# Model Architecture

❑ **Three architectures:**

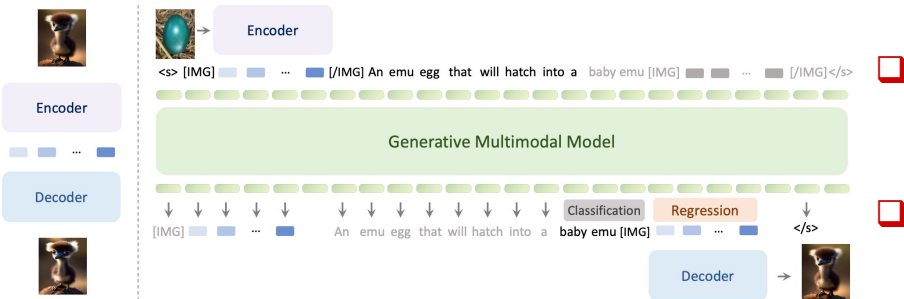**(a) language-centered method; (b) image-centered method; (c) unified method**

Wu, S., Fei, H., Qu, L., Ji, W. and Chua, T.S., 2023. Next-gpt: Any-to-any multimodal llm. ICMLR 2024.

Rust, P., Lotz, J.F., Bugliarello, E., Salesky, E., de Lhoneux, M. and Elliott, D., 2023, September. Language Modelling with Pixels. ICLR 2023.

Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sagnak Ta¸sırlar. Introducing our multimodal models: fuyu-8b, 2023. https://www.adept.ai/blog/fuyu-8b.

# In-Context Learning



- **Each image in the multimodal sequence is tokenized into embeddings via a visual encoder, and then interleaved with text tokens for autoregressive modeling.**

- **Leveraging few-shot Prompting for diverse reasoning tasks**

Sun, Q., Cui, Y., Zhang, X., Zhang, F., Yu, Q., Luo, Z., Wang, Y., Rao, Y., Liu, J., Huang, T. and Wang, X. Generative multimodal models are in-context learners. CVPR 2024.

# Evolution of Multimodal Reasoning

❑ **From task-specific to centralized paradigms**



VITRON: A Unified Pixel-level Vision LLM for Understanding, Generating, Segmenting, Editing. https://vitron-llm.github.io/

# Evolution of Multimodal Reasoning

❑ **From (implicit) single-step prediction to (explicit) multi-step reasoning**



（a）An example of ScienceQA.



（b）An example of CoCo-MMRD.

❑ **Improved Interpretability**: offer an interpretable glimpse into the decision-making process

❑ **Improved Controllability**: interfere the reasoning process, e.g., adding complementary information, verifying and correcting mistakes

❑ **Improved Flexibility**: allow interactive communications between different models



Wei, J., Tan, C., Gao, Z., Sun, L., Li, S., Yu, B., Guo, R. and Li, S.Z., 2023. Enhancing Human-like Multi-Modal Reasoning: A New Challenging Dataset and Comprehensive Framework. arXiv preprint arXiv:2307.12626.

# 2

# Multimodal Chain-of-Thought Reasoning

# Multimodal Chain-of-Thought Reasoning

❑ **Think step by step, formulate intermediate steps before deriving an answer**

❑ **Paradigm shift of task format**

- **Standard Format: <input → output>**

- **CoT Format: <input → rationale → output>**

Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G. and Smola, A. Multimodal chain-of-thought reasoning in language models. TMLR 2024.

# Multimodal Chain-of-Thought Reasoning

❑ **Input: Various modalities such as text, image, and graph are incorporated into the model's input**

❑ **Output: Multimodalities, including text and image, are generated in the model's output**



Video-of-Thought (VoT)

Zhang, Z., Yao, Y., Zhang, A., Tang, X., Ma, X., He, Z., Wang, Y., Gerstein, M., Wang, R., Liu, G. and Zhao, H., 2023. Igniting Language Intelligence: The Hitchhiker's Guide From Chain-of-Thought Reasoning to Language Agents. arXiv preprint arXiv:2311.11797.

Fei, H., Wu, S., Ji, W., Zhang, H., Zhang, H., Lee, M., Hsu, W., 2024. Video-of-Thought: Step-by-Step Video Reasoning from Perception to Cognition. ICML 2024.

# The Role of (Multimodal) Chain-of-Thought

❑ **Role 1: Introducing more reliable input results in more convincing reasoning process**

❑ **Case studies: 50 error cases**

- **Imperfect training data: when the vision input is missing**
- **Generate hallucinated rationales that mislead the answer inference (64%)**

### Problem

**Question:** Will these magnets attract or repel each other?
**Context:** Two magnets are placed as shown. Hint: Magnets that attract pull together. Magnets that repel push apart.

Vision: S N S N

**Options:** (A) attract   (B) repel

**Gold Rationale:** Magnets can pull or push on each other without touching. When magnets attract, they pull together. When magnets repel, they push apart. Whether a magnet attracts or repels other magnets depends on the positions of its poles, or ends. Every magnet has two poles, called north and south. Here are some examples of magnets. The north pole of each magnet is marked N, and the south pole is marked S. If different poles are closest to each other, the magnets attract. The magnets in the pair below attract. If the same poles are closest to each other, the magnets repel. The magnets in both pairs below repel. Will these magnets attract or repel? To find out, look at which poles are closest to each other. The north pole of one magnet is closest to the south pole of the other magnet. Poles that are different attract. So, these magnets will attract each other.
**Answer:** The answer is (A).

### Baseline

**Generated Rationale:** Magnets can pull or push on each other without touching. When magnets attract, they pull together. When magnets repel, they push apart. Whether a magnet attracts or repels other magnets depends on the positions of its poles, or ends. Every magnet has two poles, called north and south. Here are some examples of magnets. The north pole of each magnet is marked N, and the south pole is marked S. If different poles are closest to each other, the magnets attract. The magnets in the pair below attract. If the same poles are closest to each other, the magnets repel. The magnets in both pairs below repel. Will these magnets attract or repel? To find out, look at which poles are closest to each other. The south pole of one magnet is closest to the south pole of the other magnet. Poles that are the same repel. So, these magnets will repel each other.
**Answer:** The answer is (B).

Others (36%)   Hallucination (64%)

Resolved (62.5%)   Unresolved (37.5%)

(a) ratio of hallucination mistakes      (b) correction rate w/ vision features

*Table 3.* Two-stage setting of (i) rationale generation (RougeL) and (ii) answer inference (Accuracy).

| Method | (i) QCM→ R | (ii) QCMR→ A |
|---|---|---|
| Two-Stage Framework | 91.76 | 70.53 |
| w/ Captions | 91.85 | 71.12 |
| w/ Vision Features | 96.97 | 84.91 |

Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G. and Smola, A. Multimodal chain-of-thought reasoning in language models. TMLR 2024.

# The Role of (Multimodal) Chain-of-Thought

❑ **Role 1: Introducing more reliable input results in more convincing reasoning process**

  ● **More accurate perception, less hallucinations during the reasoning process**

| Will these magnets attract or repel each other? | | |
|---|---|---|
| Will these magnets attract or repel each other? | **Baseline** | The **south pole** of one magnet is closest to the **south pole** of the other magnet. |
| | $X_1 = X_2 \rightarrow \{Y_1, Y_2\}$ | The **north pole** of one magnet is closest to the **south pole** of the other magnet. |

Will these magnets attract or repel each other?

N S | S N

**+ Vision**

$X_1 \rightarrow Y_1$
$X_2 \rightarrow Y_2$

Will these magnets attract or repel each other?

S N | S N

The **south pole** of one magnet is closest to the **south pole** of the other magnet.

The **north pole** of one magnet is closest to the **south pole** of the other magnet.

Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G. and Smola, A. Multimodal chain-of-thought reasoning in language models. TMLR 2024.

# The Role of (Multimodal) Chain-of-Thought

❑ **Role 2: Breaking complex problems into smaller, manageable sub-problems**

  ● **identifies atomic knowledge components essential for reasoning processes**

  ● **integrates their relationships, thereby constructing coherent reasoning steps**



CoT reasoning is useful because

1) direct prediction is inaccurate for some inferences because the relevant variables are rarely seen together in training

2) chain-of-thought reasoning improves estimation by incrementally chaining local statistical dependencies that are observed frequently in training

Prystawski, B., Li, M., & Goodman, N. Why think step by step? Reasoning emerges from the locality of experience. NeurIPS 2024.

# The Role of (Multimodal) Chain-of-Thought

☐ **Role 3: Available for stepwise knowledge update and self-correction (w/ external feedback)**

☐ **Chain-of-Thought w/ Tools**

- **Plan-Solve, augment and correct intermediate steps with tools**

- **Overcome the Intrinsic ability deficiency of LLMs such as calculation, searching**

### Question Answering

**Question**: Who won the 2016 Russian national silver medal with another Russian ice dancer born 29 April 1995?
**Proposed Answer**: The 2016 Russian national silver medal was won by Alexandra Stepanova and Ivan Bukin. Alexandra Stepanova was born 29 April 1995. Ivan Bukin was born 10 April 1993. So the answer is: Ivan Bukin. ✗

What's the problem with the above answer?
**1. Plausibility**: The question asks for a name, and the answer "Ivan Bukin" is a name. So it's plausible.

**2. Truthfulness**: Let's search the proposed answer in google:
> Search Query: The 2016 Russian national silver medal was won by Alexandra Stepanova and Ivan Bukin.
> Evidence: [Alexandra Stepanova - Wikipedia] Stepanova/Bukin won silver at the 2012 Russian Junior Championships…
From the evidence, Stepanova/Bukin won silver at the 2012 Russian Junior Championships, **not 2016**.

Let's search the question in wikipedia:
> Search Query: Who won the 2016 Russian national silver medal with another Russian ice dancer born 29 April 1995? site: wikipedia.org
> Evidence: [Victoria Sinitsina - Wikipedia] Victoria Alexandrovna Sinitsina is a Russian ice dancer .

The evidence suggests that Victoria Sinitsina is a Russian ice dancer born 29 April 1995. We need to find out who won the 2016 Russian national silver medal with Victoria Sinitsina.

> Search Query: Who won the 2016 Russian national silver medal with Victoria Sinitsina?
> Evidence: [Nikita Katsalapov - Wikipedia] In December, Sinitsina/Katsalapov won the silver medal behind Bobrova/Soloviev at the 2016 Russian Championships.

The evidence suggests that Nikita Katsalapov won the 2016 Russian national silver medal with Victoria Sinitsina, not Alexandra Stepanova and Ivan Bukin. The answer is **incorrect**.

**Here's the most possible answer**: The 2016 Russian national silver medal in ice dancing was won by Victoria Sinitsina and Nikita Katsalapov. Victoria Sinitsina was born on April 29, 1995. So the answer is: Nikita Katsalapov. ✓

### Program Synthesis

**Question**: Ann, Bill, Cate, and Dale each buy personal pan pizzas cut into 4 pieces. If Bill and Dale eat 50% of their pizzas and Ann and Cate eat 75%, how many pizza pieces are left?
**Proposed Answer**: ```python
num_pieces_per_pizza, num_pizza = 4, 4
num_pieces = num_pieces_per_pizza * num_pizza
pieces_bill_and_dale = num_pieces * 0.5
pieces_ann_and_cate = num_pieces * 0.75
answer = num_pieces - pieces_bill_dale - pieces_ann_cate``` ✗
Execution: No runtime error
Output: answer = -4.0

What's the problem with the above code?
**1. Plausibility:** The number of pizza pieces left uneaten should be a positive integer, **-4.0 < 0**, so the answer is **not reasonable**.

**2. Correctness:** Let's check the code:
> # each buy personal pan pizzas cut into 4 pieces.
> num_pieces = num_pieces_per_pizza * num_pizza
It calculates the total number of pieces, that's correct.

> pieces_bill_dale = num_pieces * 0.5
> pieces_ann_cate = num_pieces * 0.75
Each person only eats their own personal pan pizza, `num_pieces * 0.5` means 50% of the total pieces, this is **wrong**.

**Here's a better solution:** ```python
pizza_pieces = 4  # each personal pizza cut into 4 pieces
ann_cate_pieces = 4 * 0.75 * 2 # Ann, Cate eats 75%
bill_dale_pieces = 4 * 0.5 * 2 # Bill, Dale eats 50%
total_eaten = ann_cate_pieces + bill_dale_pieces
answer = pizza_pieces * 4 - total_eaten # pieces left``` ✓

### Toxicity Reduction

**Previous toxic text:** … last year before hanging himself. ✗
The text has 53% toxicity of threat. To improve content safety while ensuring fluency, here's a **non-toxic continuation** after '\':

Jason Conroy, 19, strangled fellow resident Melissa Mathieson at Alexandra House in Bristol \ in November 11 2014. It serves as a reminder of the importance of mental health support to… ✓

Gou, Zhibin, et al. "Critic: Large language models can self-correct with tool-interactive critiquing. ICLR 2024.

# 3 Towards Multimodal LLM Agents

# Towards Multimodal LLM Agents

❑ **From content-based reasoning to behavior control (w/ multimodalities)**

❑ *"Those who know but do not act simply do not yet know"*



Brain in a Vat

Ma, Y., Zhang, C. and Zhu, S.C., 2023. Brain in a vat: On missing pieces towards artificial general intelligence in large language models. arXiv preprint arXiv:2307.03762.
Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E. and Zheng, R., 2023. The rise and potential of large language model based agents: A survey. arXiv preprint arXiv:2309.07864.

# Towards Multimodal LLM Agents

❑ **(M)LLM Agents:** follow language instructions and execute actions in environments, possibly use tools

❑ **Features:** General, Autonomous, Adaptive, Evolutionary, Socialized

# Towards Multimodal LLM Agents



**Control: OS and Applications**



**Research: Organic Synthesis**



**Programming: Code Generation**



**Control: Embodied Systems**



**Research: Medical Assistance**



**Interaction: Multi-Agent Collaboration**

Ma, Y., Zhang, C. and Zhu, S.C., 2023. Brain in a vat: On missing pieces towards artificial general intelligence in large language models. arXiv preprint arXiv:2307.03762.
Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E. and Zheng, R., 2023. The rise and potential of large language model based agents: A survey. arXiv preprint arXiv:2309.07864.

# Taxonomy of (M)LLM Agents

## Autonomous Agents

**ADEPT** **Action Transformer**
https://www.adept.ai/blog/act-1

**Google** **AITW**
https://github.com/google-research/google-research/tree/master/android_in_the_wild

**WebArena**
https://webarena.dev

**Auto-UI**
https://github.com/cooelf/Auto-UI

## Communicative Agents

**CAMEL**
https://github.com/camel-ai/camel

**Generative Agents**
https://github.com/joonspk-research/generative_agents

**VOYAGER**
https://voyager.minedojo.org/

**ChatDev**
https://github.com/OpenBMB/ChatDev

More: AutoGPT, BabyAGI, Meta-GPT, AgentGPT

## Autonomous Agents: mainly task automation

**Mobile Device Automation**



Meta-GUI

**Webpage Automation**



WebArena

**Application Automation**



ACT-1

Sun, Liangtai, et al. "META-GUI: Towards Multi-modal Conversational Agents on Mobile GUI." *EMNLP 2022.*
Zhou, Shuyan, et al. "Webarena: A realistic web environment for building autonomous agents." *arXiv preprint arXiv:2307.13854* (2023).
*https://www.adept.ai/blog/act-1*

# Taxonomy of (M)LLM Agents

## Communicative Agents: personalized, socialized, interactive

### Agents-Agents

### Agents-Human



Park, Joon Sung, et al. "Generative agents: Interactive simulacra of human behavior." *arXiv preprint arXiv:2304.03442* (2023).
Lin, Jessy, et al. "Decision-Oriented Dialogue for Human-AI Collaboration." *arXiv preprint arXiv:2305.20076* (2023).

# Technological Paradigm

**Task Instruction**

**Environment**

| OS | APP |
| Webpage | Virtual Env. |

Interaction

**Tool**

| API Interface | Physical Device |
| Rule Set | Interpreter |

Planning / Problem Decomposition

**(M)LLM**

Plan

Decision

Control

Memory (long/short)

State

Action

Execute / Call

Decision Making

**Foundation**
- ❏ Multimodalities
- ❏ Long-context Modeling

**Workflow**
- ❏ Perception
- ❏ Planning & Decision Making
- ❏ Action (w/ Tool Use)
- ❏ Interaction
- ❏ Memory
- ❏ Multi-Agent Collaboration

Act

Obs

# GUI Agents

❑ **Auto-GUI：Multimodal Autonomous Agents for GUI control**

● **assist users in completing tasks in distinct environments such as operation systems, specific applications, and web browsers**

● **Imitate human clicking, scrolling, and typing actions, and operate directly with the GUI**



Goal: turn off javascript in the chrome app

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.
Xinbei Ma, Zhuosheng Zhang, Hai Zhao. Comprehensive Cognitive LLM Agent for Smartphone GUI Automation. Findings of ACL 2024.
https://machinelearning.apple.com/research/ferret..

# Auto-UI

- **Multimodal Agent: BLIP2 + FLAN-Alpaca**

- **Chain-of-Action: a series of intermediate previous action histories and future action plans**



**Goal:** Look up the best rated coffee maker on Lowe's    $X_{goal}$

**Chain of Previous Action Histories:**
action_type: type, touch_point: [-1.0, -1.0], lift_point: [-1.0, -1.0], typed_text: "best rated coffee maker"
action_type: dual_point, touch_point: [0.2, 0.5], lift_point: [0.8, 0.5], typed_text: ""    $X_{history}$

$X_{language}$

Chain of Actions

$X_{screen}$

Image Encoder

Language Encoder

Projection

Self Attention

Feedforward

Decoder

Chain of Future Action Plans

**Action Plan:**
[DUAL_POINT,
STATUS_TASK_COMPLETE]    $Y_{plan}$

Current Action Prediction

**Action Decision:**
action_type: [DUAL_POINT],
touch_point: [0.5595, 0.6261],
lift_point: [0.5595, 0.6261], typed_text: ""    $Y_{action}$

Screen

Action

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

28

# Results

- A **unified multimodal model** out of *first principles thinking* can serve as a strong autonomous agent

  - can be adapted to **different scenarios** without the need to train specific models for each task

  - does not need additional annotations (screen parsing) and is **easy to use**

- Coverage: 30K unique instructions, 350+ Apps and websites

- **Action Type Accuracy: 90%+, Action Success Rate: 74%+**

| Model | Unified | w/o Anno. | Overall | General | Install | GoogleApps | Single | WebShopping |
|-------|---------|-----------|---------|---------|---------|------------|--------|-------------|
| BC-single | ✗ | ✗ | 68.7 | - | - | - | - | - |
| BC-history | ✗ | ✗ | <u>73.1</u> | <u>63.7</u> | <u>77.5</u> | <u>75.7</u> | <u>80.3</u> | <u>68.5</u> |
| PaLM 2-CoT | ✓ | ✗ | 39.6 | - | - | - | - | - |
| ChatGPT-CoT | ✓ | ✗ | 7.72 | 5.93 | 4.38 | 10.47 | 9.39 | 8.42 |
| Fine-tuned Llama 2 | ✗ | ✗ | 28.40 | 28.56 | 35.18 | 30.99 | 27.35 | 19.92 |
| Auto-UI$_{separate}$ | ✗ | ✓ | 74.07 | 65.94 | **77.62** | **76.45** | 81.39 | 69.72 |
| Auto-UI$_{unified}$ | ✓ | ✓ | **74.27** | **68.24** | 76.89 | 71.37 | **84.58** | **70.26** |

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

# Insights

❑ The bottleneck seems to be the **multimodal perception**, misleading the reasoning process

   ● GUI involves comprehensive elements (interleaved, icons, texts, boxes)

   ● Changing vision encoders influences the performance dramatically

❑ Scaling does not always improve performance

| Model | Overall | General | Install | GoogleApps | Single | WebShopping |
|---|---|---|---|---|---|---|
| Auto-UI on CLIP | 71.84 | 66.28 | 74.40 | 69.71 | 81.60 | 67.23 |
| Auto-UI on BLIP-2 | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |
| Auto-UI on Vanilla-T5$_{large}$ | 72.98 | 66.61 | 75.40 | 70.86 | 83.47 | 68.54 |
| Auto-UI on FLAN-T5$_{large}$ | 73.36 | 67.59 | 76.35 | 70.71 | 83.01 | 69.12 |
| Auto-UI on FLAN-Alpaca$_{large}$ | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |
| Auto-UI on FLAN-Alpaca$_{small}$ | 71.38 | 65.26 | 74.90 | 68.70 | 81.20 | 66.83 |
| Auto-UI on FLAN-Alpaca$_{base}$ | 72.84 | 66.97 | 75.93 | 70.29 | 82.56 | 68.46 |
| Auto-UI on FLAN-Alpaca$_{large}$ | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |

# Insights

❑ **Category Accuracy:** the major challenges lie within the click region and scroll direction predictions

- The model tends to click a wrong place or scroll in a wrong direction

❑ Challenge in "really" understanding the GUI layouts, e.g., relationship between GUI elements



Legend: Click (67.4%), Scroll (82.0%), Action Type (90.1%), Typed Text (93.1%)

Categories: General, Install, GoogleApps, Single, WebShopping. Y-axis: Accuracy (%).

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.
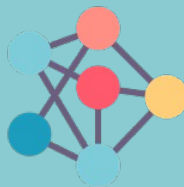
**4** Challenges

# Challenges

❑ Multimodal reasoning drives smart MLLMs

● More broader scenarios (physical and virtual worlds)

● More comprehensive scenarios (evolutionary, interactive)

**Evolutionary Reasoning**
- **Active explore and evolve from environments**
- **Learn from (un)successful attempts**

**Interactive Reasoning**
- **Human-in-the-loop interference**
- **Error identification and correction abilities**

**Reasoning Alignment**
- **Align both content safety, and behavior safety**
- **Decide the action trajectory with foresights**

33

# Summary

- **Basics of Multimodal Reasoning**
  - **Concept: derive high-level conclusions from multiple modalities, possibly via multiple logical steps based on atomic evidences**
  - **Developments: (a) From task-specific to centralized paradigms; (b) From single-step prediction to multi-step reasoning**
  - **Popular Approaches: (a) In-Context Learning: (b) Multimodal Chain-of-Thought**

- **Multimodal Chain-of-Thought Reasoning**
  - **Paradigm Shift: From "<input → output>" to <input → rationale → output>**
  - **Role 1: Introducing more reliable input results in more convincing reasoning process**
  - **Role 2: Breaking complex problems into smaller, manageable sub-problems**
  - **Role 3: Available for stepwise knowledge update and self-correction (w/ external feedback)**

- **Towards Multimodal LLM Agents**
  - **Taxonomy: Autonomous Agents and Communicative Agents**
  - **Technical Components: Foundation (multimodality & long-context modeling); (b) Workflow (plan, act, memory, feedback)**

- **Challenges**
  - **Evolutionary Reasoning, Interactive Reasoning, Reasoning Alignment**

# Thanks!

**Any questions?**
You can find me at:

✛ zhangzs@sjtu.edu.cn