
From Tools to Agents: Key Technologies for Language-Driven Action Agents

Prof. Chia-Hui Chang (張嘉惠)

National Central University, Taiwan

Convenor for the Intelligence Computing Discipline @ NSTC

Co-Chair for AI Center of Excellence @ NSTC

Jan. 20, 2025





About Me

- Dept. of CSIE, NCU
- Office of Institutional Research, NCU
- Convenor for the Intelligence Computing discipline @ NSTC
- Co-Chair for AI Center of Excellence @ NSTC



<https://www.aclclp.org.tw/>



The Association for Computational Linguistics and Chinese Language Processing 第二十三卷第二期

E-Mail : aclclp@hp.cs.sinica.edu.tw

地 址：台北市研究院路二段128號中研院資訊所

Website : <http://www.aclclp.org.tw>

電 話：(02)2788-3099 ext:1582

發行人：許開廉
主編：王新民
執行編輯：黃琪

傳 真：(02)2788-1638
郵 箱：18066251



中華民國人工智慧學會
Taiwanese Association for Artificial Intelligence

<https://www.taai.org.tw/>

Outline

- Background
 - 4 Phases of AI
- What are AI Agents?
 - Example: GUI Agents, Robots
 - Agents for Software Development
- Course Administration, Virtual TA by EduACT
- Your turn



PHYSICAL AI
SELF-DRIVING CARS
GENERAL ROBOTICS



AGENTIC AI
CODING ASSISTANT
CUSTOMER SERVICE
PATIENT CARE

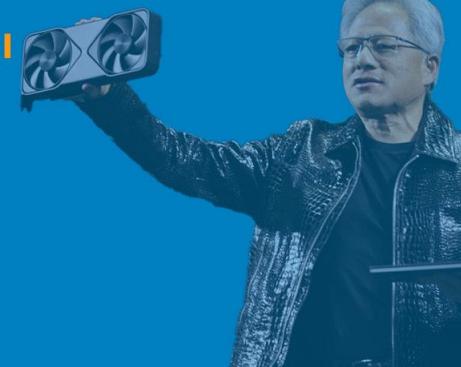


GENERATIVE AI
DIGITAL MARKETING
CONTENT CREATION



PERCEPTION AI
SPEECH RECOGNITION
DEEP RECSYS
MEDICAL IMAGING

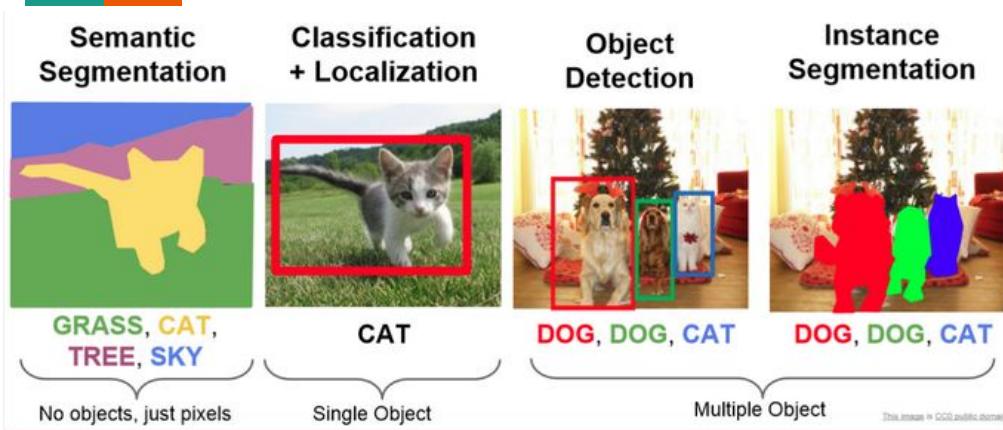
2012 ALEXNET



The Evolution of AI: From AlexNet to Reality

NVIDIA CEO Jensen Huang
Keynote at CES 2025

Vision: Perception vs. Generative



"A brain riding a rocketship heading towards the moon."

text-to-image model



Natural Language Understanding

Elon Musk PERSON apparently wasn't aware that his company SpaceX had a Facebook ORG page. The SpaceX and Tesla PRODUCT CEO has responded to a comment on Twitter GPE calling for him to take down the SpaceX, Tesla and Elon Musk ORG official pages in support of the #deletefacebook movement by first ORDINAL acknowledging he didn't know one existed, and then following up with promises that he would indeed take them down.

He's done just that, as the SpaceX NORP Facebook page is now gone, after having been live earlier today DATE (as you can see from the screenshot included taken at around 12:10 PM ET) TIME .

- Named Entity Recognition
- Relation Extraction
- Event Extraction
- Sentiment Analysis
- Part of Speech Tagging
- Chunking
- Dependency parsing

Natural Language Generation

Conversation extract:

User I want to review a film

LUCID Sure, I'd be happy to help with that. Can you please tell me the name of the film you'd like to review?

User Bobby, stop playing with your food and eat your dinner!

LUCID It seems like you're busy at the moment. Let's continue this conversation after dinner.

User Sorry about that, my son was misbehaving. The film I want to review is The Dark Knight

LUCID Great, you would like to review "The Dark Knight". Is that correct?

User Yes, that's correct. My review is: I found the storyline of this movie a bit hard to follow.

Ferret



Fuyu-8B



Gemini



GPT-4v



Chameleon



MM1



Video-LLaMA



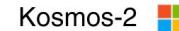
Claude 3



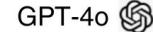
HuggingGPT



Kosmos-2



GPT-4o



LLaVA



AnyMAL



Gemini 1.5



Kosmos-1



PaLM-E



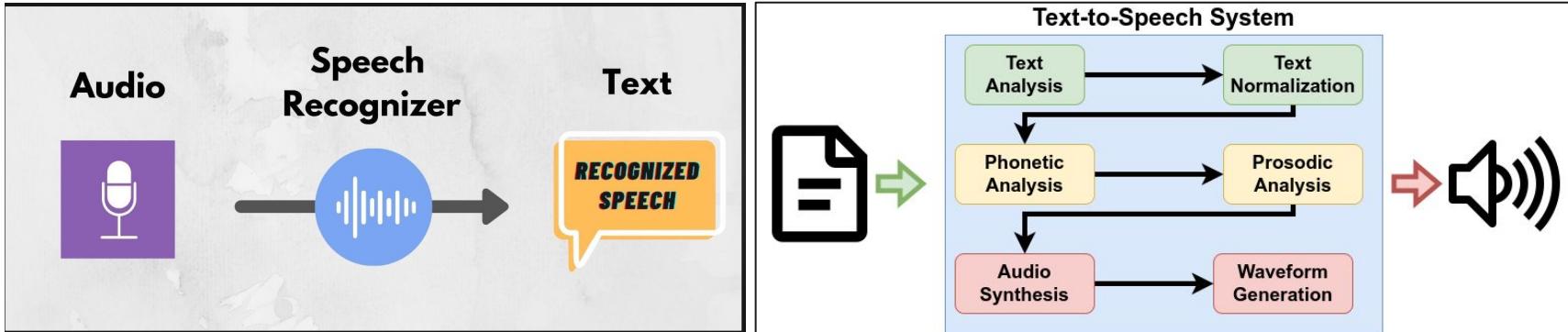
2023



2024



Speech: Recognition vs. Synthesis



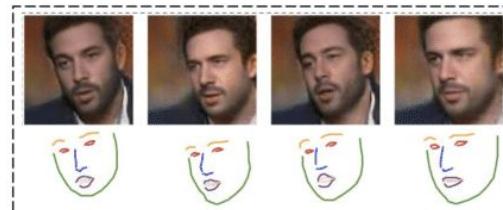
With image generation



Audio driving



(a) Audio-driven talking-head generation



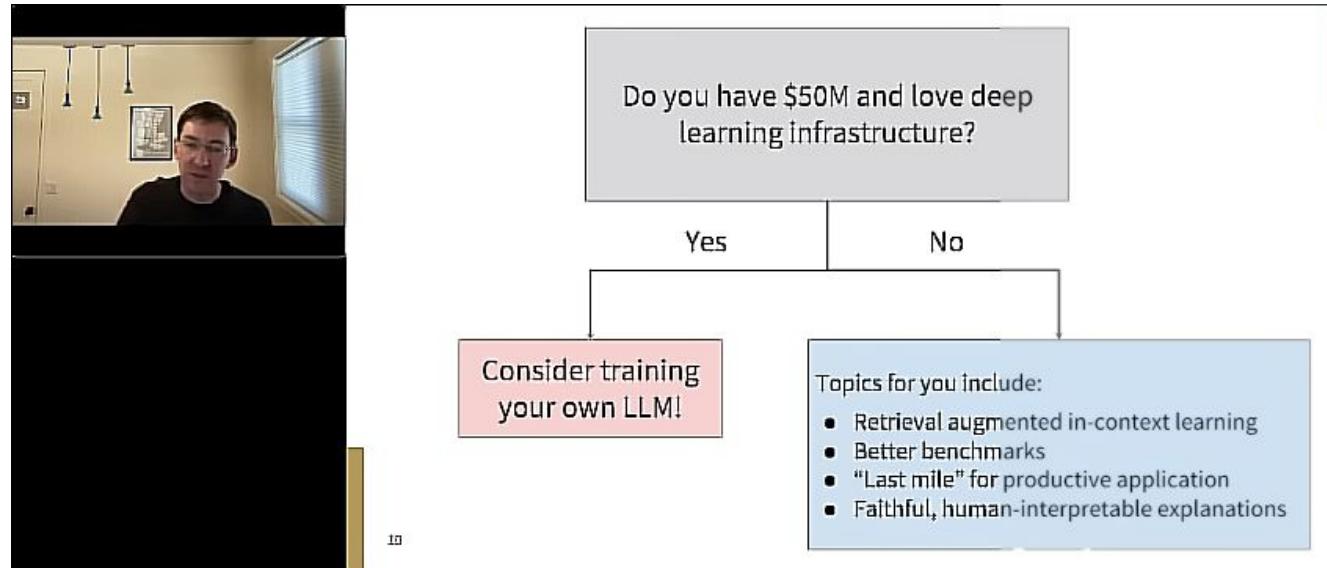
Motion driving



(b) Motion-driven talking-head generation

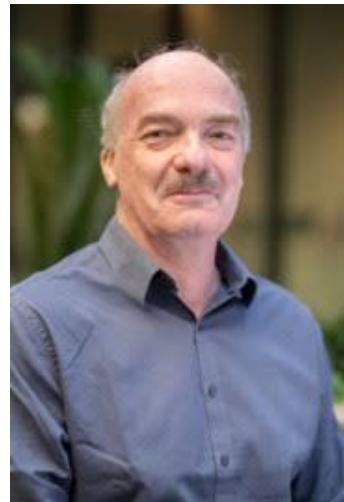
What we can contribute to NLU after GPT3?

Christopher Potts@Stanford



What are the Last Mile productive applications?

We're in a GenAI world ... let's do responsible work

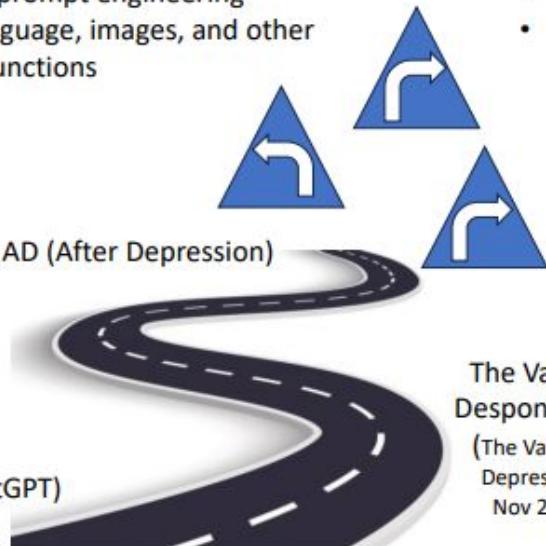


NLP engineering: Make LLMs usable

- Build smaller and cheaper LLMs
- Systematize prompt engineering
- Integrate language, images, and other media and functions

AD (After Depression)

BC (Before ChatGPT)



NLP applications: Make LLMs useful

- Tune LLMs to domains and companies for enterprise processing
- Add functionality and agency in the world
- Tailor LLMs to people to be their personal daemons/amanuenses in everyday life

NLP research: Make LLMs understandable

(or at least, be solid engineering)

- Fix the problems with LLMs
- Get explanations how LLMs do what they do
- Formalize them well enough for autonomy, assurance, and ethics



Audience Q&A

- ⓘ Presenting with animations, GIFs or speaker notes? Enable our [Chrome extension](#)

AI Agents

The Rise and Potential of Large Language Model Based Agents: A Survey

<https://arxiv.org/pdf/2309.07864>

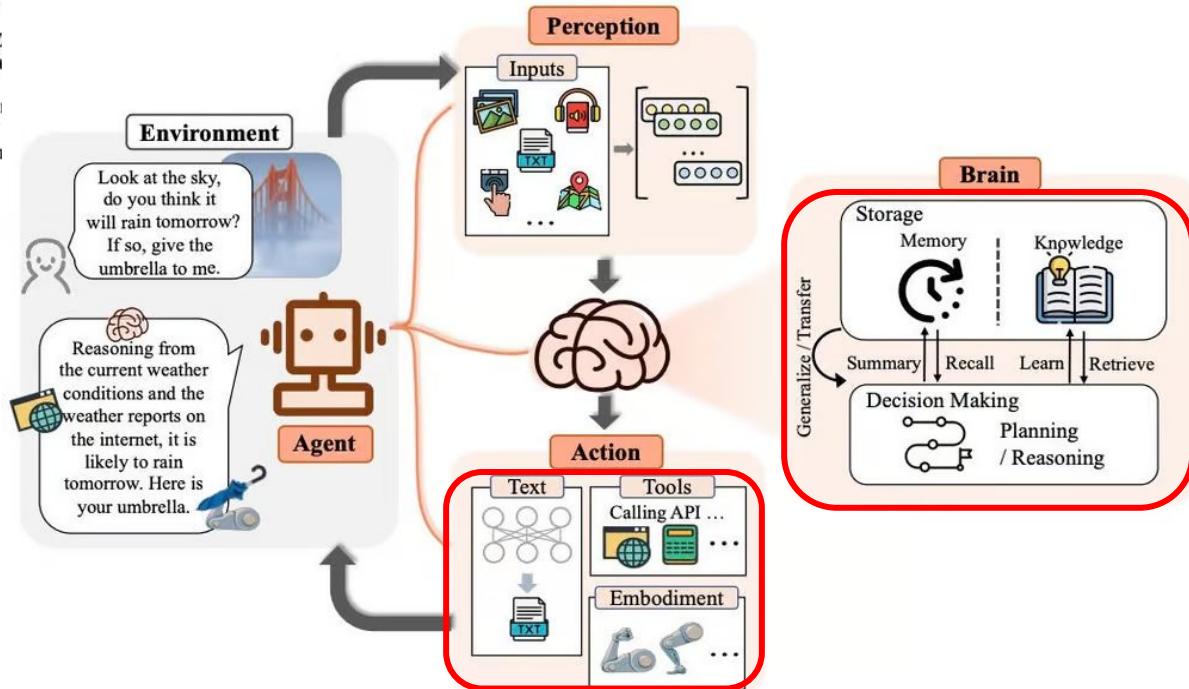
Zhiheng Xi^{*†}, Wenxiang Chen*, Xin Guo*, Wei He*, Yiwen Ding*, Boyang Hong*, Ming Zhang*, Junzhe Wang*, Senjie Jin*, Enyu Zhou*

Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zou, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhao

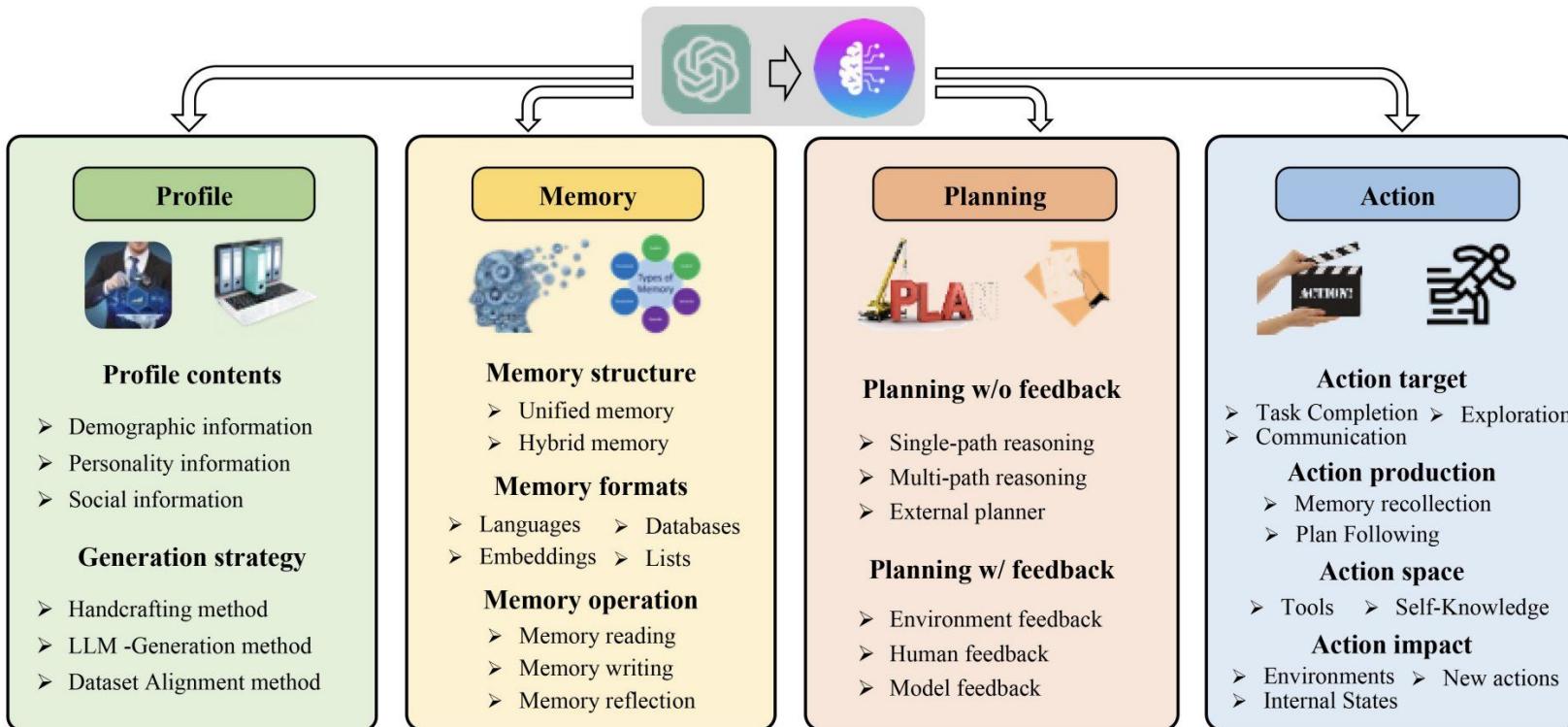
Shihang Dou, Rongxiang Weng, Wensen Chen

Qi Zhang[†], Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjin

Fudan NLP Group



Structure of an Agent



Agentic Reasoning Design Patterns

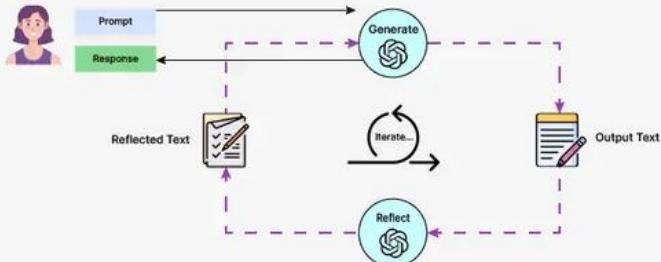
By Andrew Ng @ Sequoia AI Ascent 2024

- Planning
- Tool Use
- Reflection
- Multi-agent collaboration

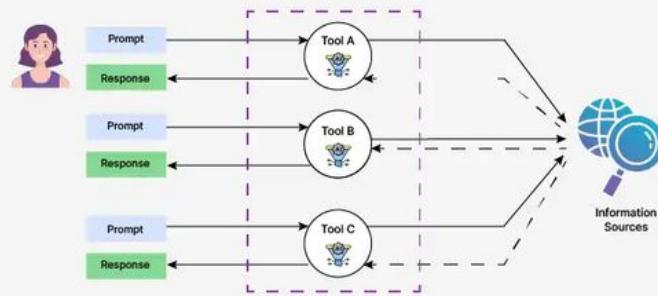


Agentic Design Patterns

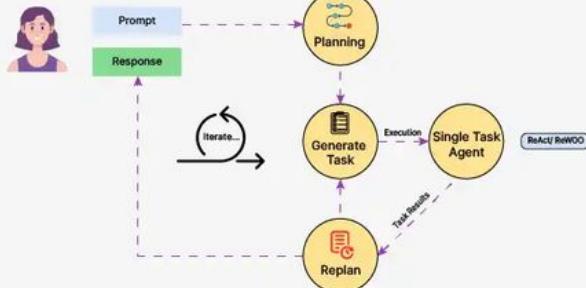
Reflection Pattern



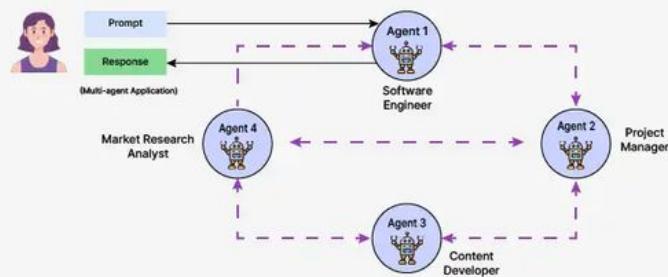
Tool Use Pattern



Planning Pattern



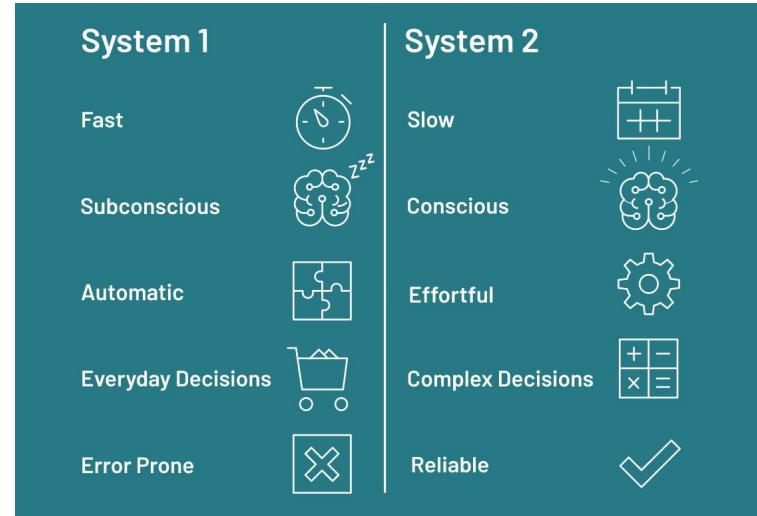
MultiAgent Pattern



Reflection Agents

Reflection is a prompting strategy used to improve the quality and success rate of agents and similar AI systems..

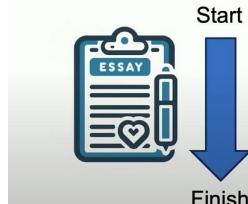
Reflection Agents



LLM-based agents

Non-agentic workflow (zero-shot):

Please type out an essay on topic X from start to finish in one go, without using backspace.



Agentic workflow:

Write an essay outline on topic X

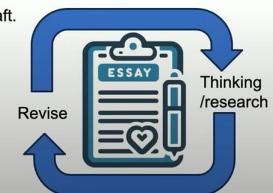
Do you need any web research?

Write a first draft.

Consider what parts need revision or more research.

Revise your draft.

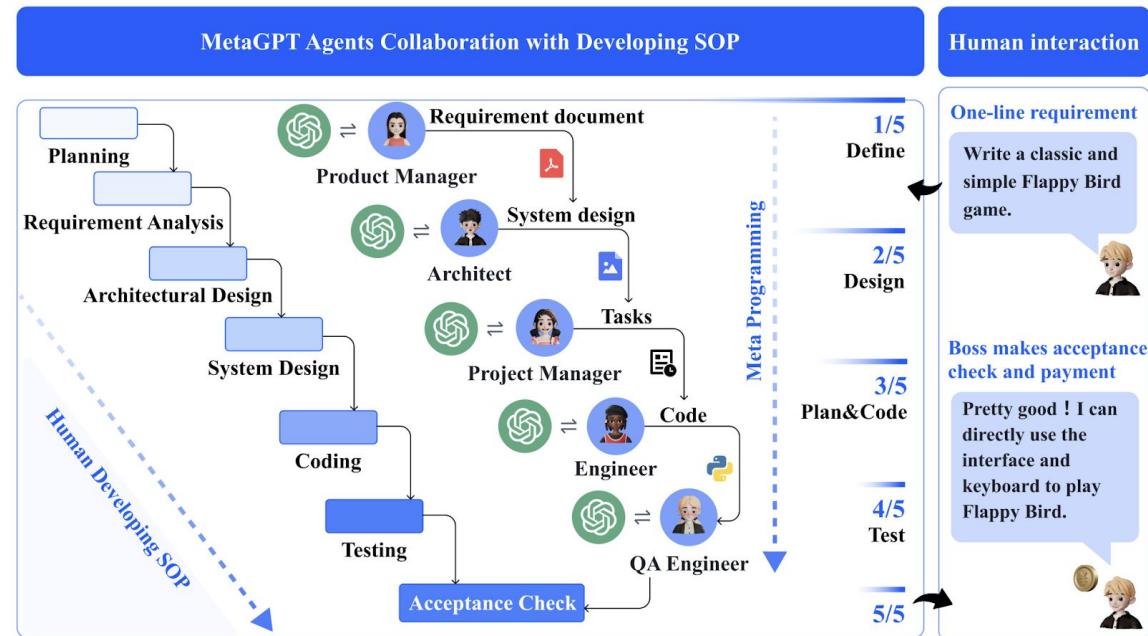
....



Multi-Agent Collaboration

MetaGPT:
AI Software Company
where LLMs act as

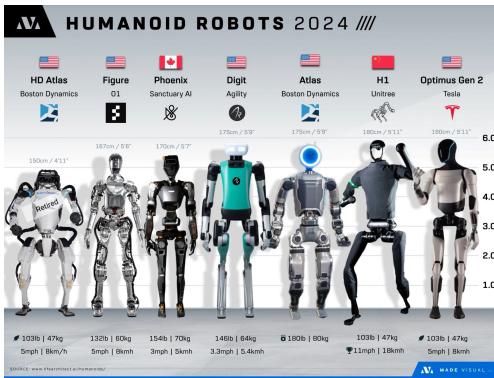
- Product Manager
- Architect
- Project Manager
- Engineer
- QA Engineer



Environment: Three Worlds

The Environment: Three Worlds

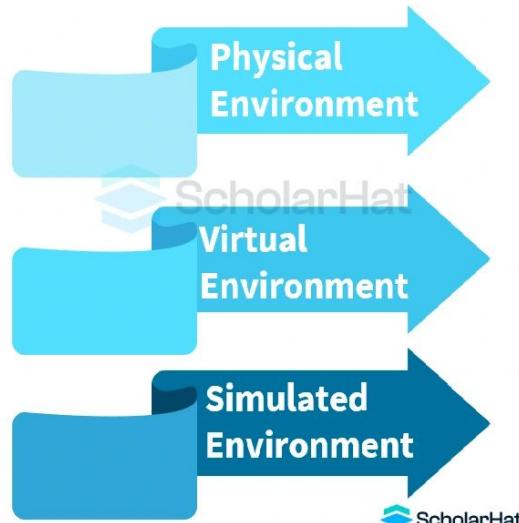
Physical world



Digital world



Types of Environment in AI



Omniverse world

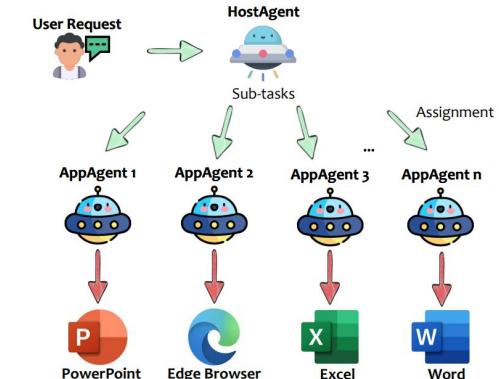


Fig. 21: The multi-agent architecture employed in UFO [17]. Figure adapted from the original paper.

Digital World



Augmenting large language models with chemistry tools

<https://arxiv.org/abs/2304.05376>

Andres M. Bran^{12*} Sam Cox^{3*} Oliver Schiltner²⁴
Carlo Baldassari¹ Andrew D. White³ Philippe Schwaller¹²

¹ Laboratory of Artificial Chemical Intelligence (LIAC), ISIC, EPFL

² National Centre of Competence in Research (NCCR) Catalysis, EPFL

³ Department of Chemical Engineering, University of Rochester

⁴ Accelerated Discovery, IBM Research – Europe

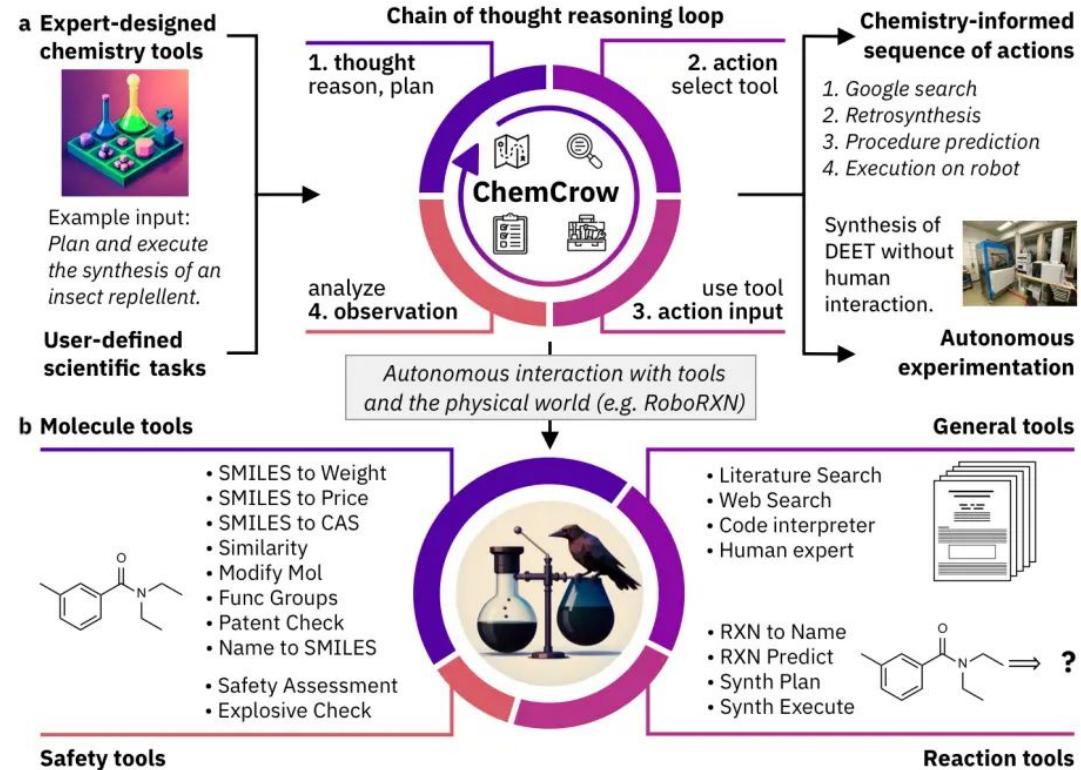
*Contributed equally.

andrew.white@rochester.edu

philippe.schwaller@epfl.ch

Require APIs to

- Web Search
- LitSearch
- Python REPL
- Name2SMILES
- SMILES2Price
- Name2CAS
- Similarity
- PatentCheck
- FuncGroups
- SMILES2Weight



Gorilla: Large Language Model Connected with **Massive APIs**

NeurIPS 2024 poster, UC Berkeley

Systems and Algorithms for
Integrating LLMs with Applications,
Tools, and Services



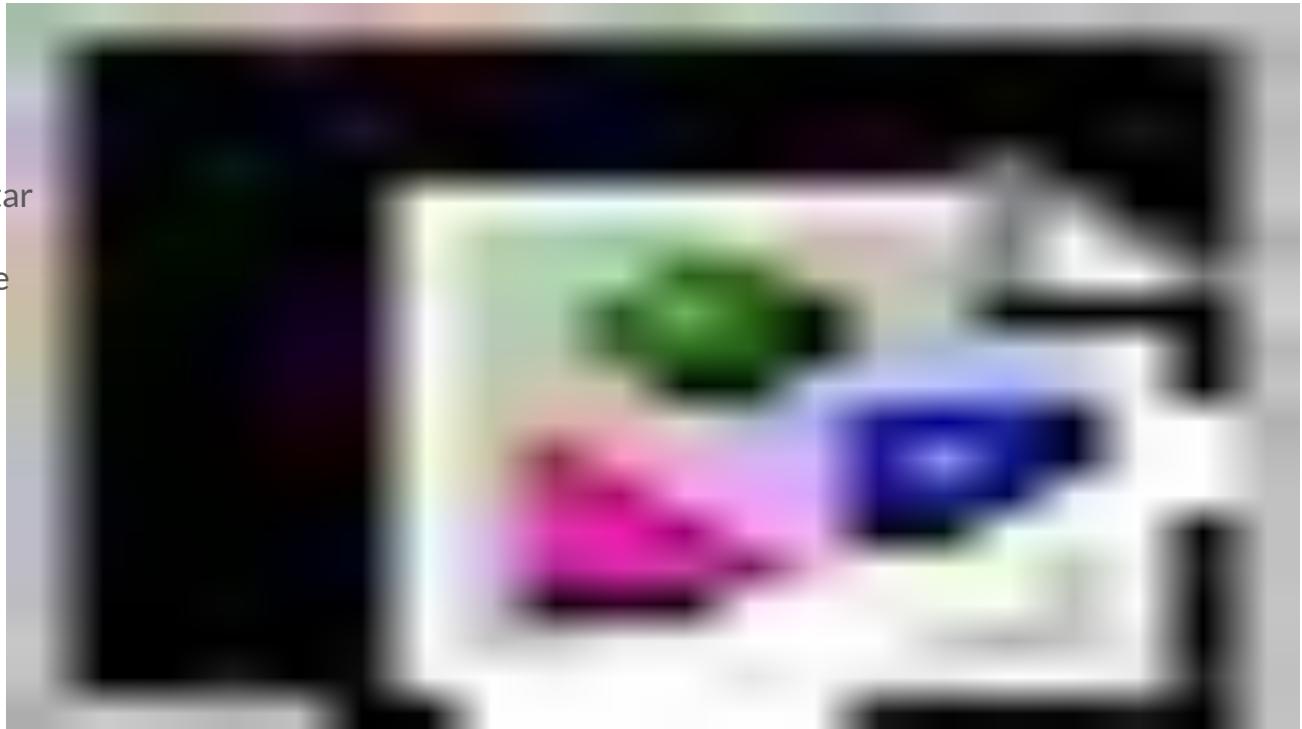
Rather have the user at the center, Gorilla enables users to interact with a wide range of services through LLMs. Gorilla is an open-source, state-of-the-art LLM that invokes API calls to interact with services!

Gorilla: Large Language Model Connected with **Massive APIs**

[Gorilla Open Functions v2](#)

1. Ask about Github Star history
2. Gorilla output valide REST API
3. Navigate to website

- [Colab](#)



Code Generation

Multiple Functions	Parallel Functions	Function Relevance Detection	Rest API
<p>User: Prompt: What is 2 + 3?</p> <p>Function: [add(int a, int b), mult(int a, int b)]</p> <p>Agent: <code>add(a=2, b=3)</code></p>	<p>User: Prompt: What is (2 + 3) and (4 + 5)?</p> <p>Function: [add(int a, int b)]</p> <p>Agent: <code>[add(a=2, b=3), add(a=4, b=5)]</code></p>	<p>User: Prompt: What is 2*3?</p> <p>Function: [add(int a, int b)]</p> <p>Agent: Error. The user asks for adding but we only have multiplication.</p>	<p>User: Prompt: get weather in Berkeley?</p> <p>Function: [requests.get(str url, str loc)]</p> <p>Agent: <code>requests.get(url=, loc = 'Berkeley')</code></p>

What if no API is available?

e.g. existing tools that are designed for humans WWW

CogAgent: A Visual Language Model for GUI Agents

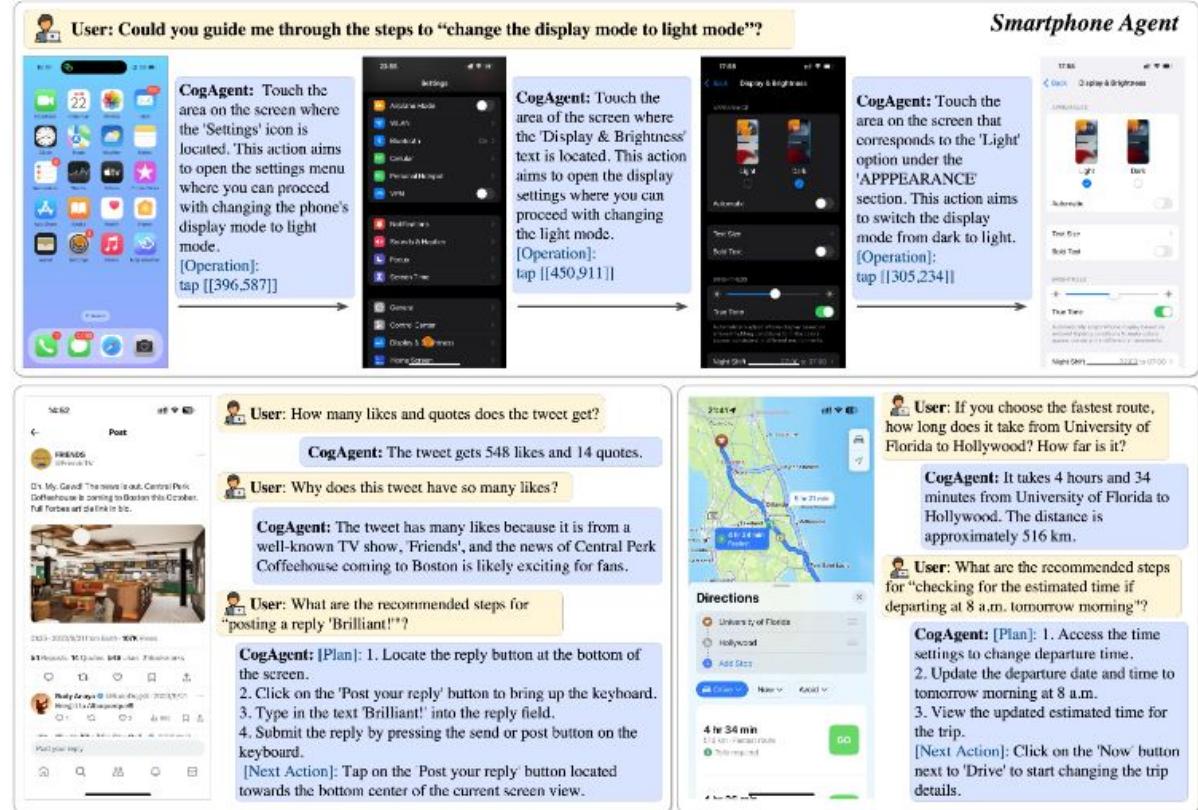
[CVPR2024]

CogAgent

- Visual Language Model based on CogVLM.
- CogAgent-18B 11B visual parameters and 7B language parameters

GUI agents for smartphones are important especially for busy persons or people with poor eyesight or elderly people

<http://36.103.203.44:7861/>



Large Language Model-Brained GUI Agents: A Survey

Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li, Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue Liu, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, Qi Zhang

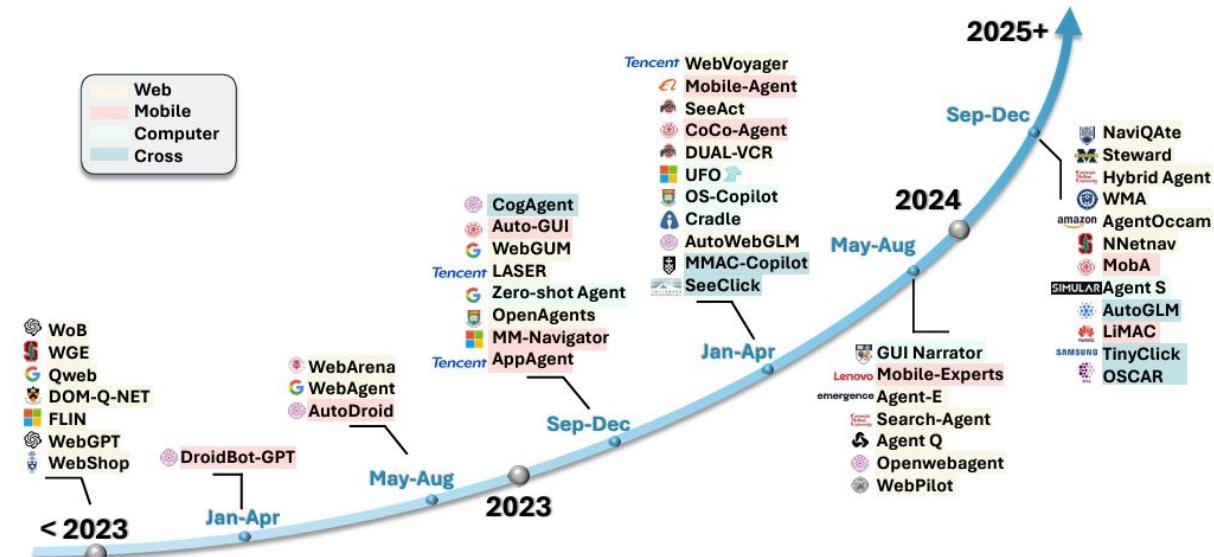
<https://arxiv.org/html/2411.18279>

From

Graphical UI

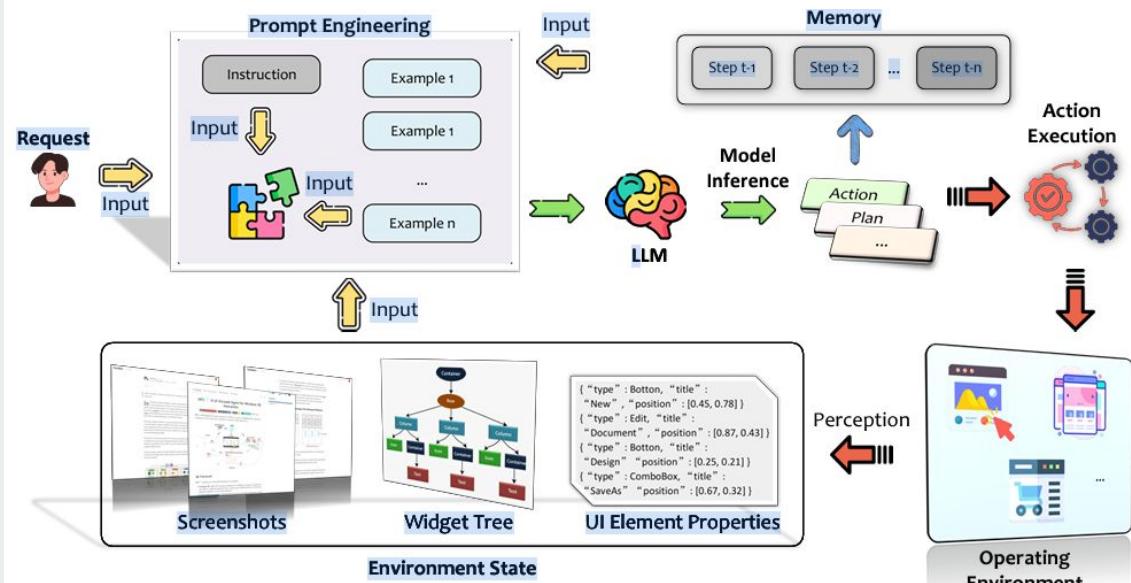
to

Conversational UI



Web Agents

From
Graphical UI
to
Conversational UI





SeeClick: Harnessing GUI Grounding for Advanced Visual GUI Agents

Kanzhi Cheng^{◊♡ *}

Qiushi Sun[♡]

Yougang Chu[◊]

Fangzhi Xu[♡]

Yantao Li[◊]

Jianbing Zhang^{◊ †}

Zhiyong Wu^{♡ †}

◊National Key Laboratory for Novel Software Technology, Nanjing University

♡Shanghai AI Laboratory

- SeeClick only relies on screenshots for task automation.

Instruction: Download the e-receipt with the last name Smith and confirmation number X123456989.

Text-based:

```
<form element_id="200">
...
<label element_id="205">Last Name:</label>
<input type="text" name="lastname" element_id="206">
...
<input type="submit" value="Get Receipt" element_id="210">
```

Simplified HTML Code

 Text-based agent's next action

Element: <element_id=206>

Action: CLICK

Selenium Code
element = driver.find_element(By.XPATH,
'//*[@@element_id="206"]')
element.click()

Vision-based:



GUI Screenshot

 SeeClick's next action

{"action": "click", "loc": [0.46, 0.62]}

Car Rentals e-Receipts

To request a receipt, please complete the fields below, or log in to your Budget profile and access your Past Rentals page.
Note: Microsoft® can be added when viewing receipt.

Country: Last Name: Confirmation/Rental Number:

United States ▾ Last Name: Confirmation/Rental Number:

Get Receipt

Please make sure Adobe Reader is required to view your receipt. If you don't have Adobe Reader, Please click here to download it.

Last Name: Confirmation/Rental Number:

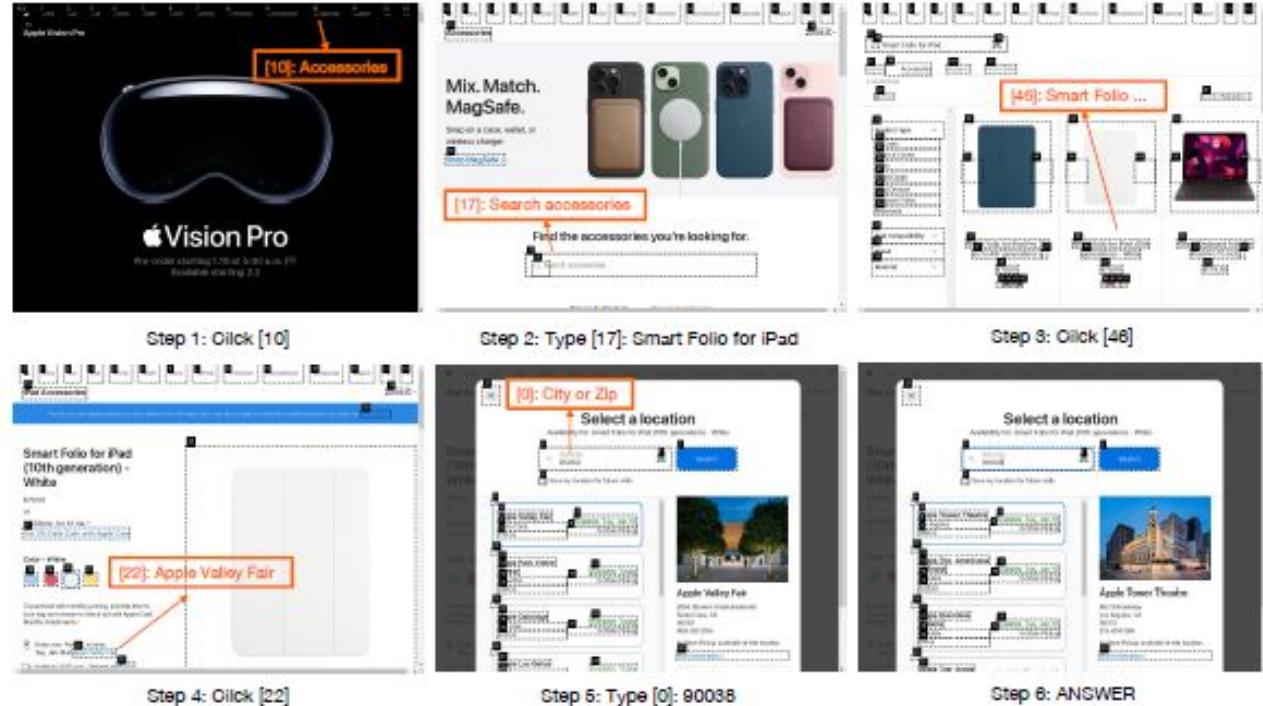
Get Receipt

WebVoyager : Building an End-to-End Web Agent with Large Multimodal Models

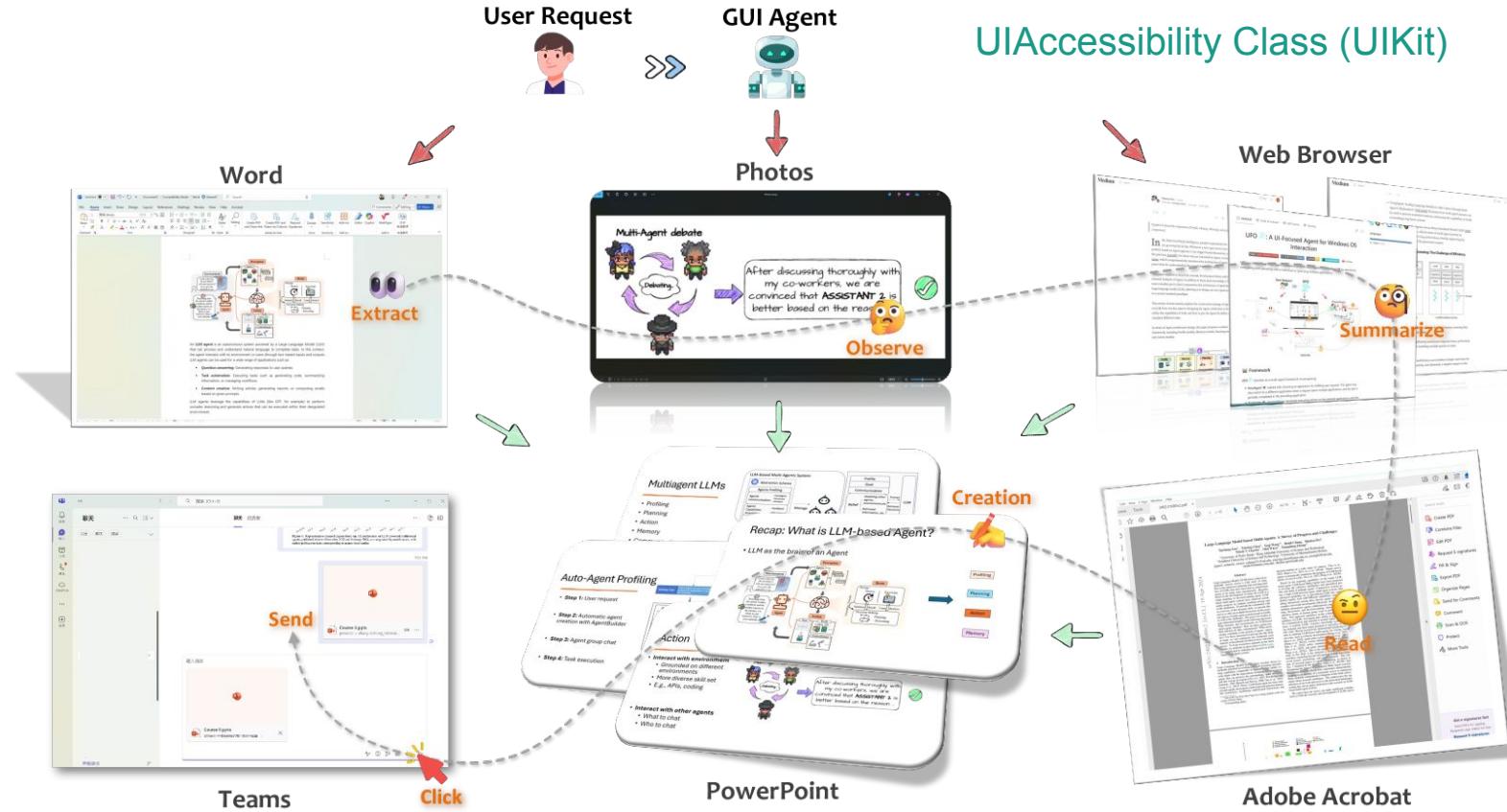
Hongliang He^{1,3*}, Wenlin Yao², Kaixin Ma², Wenhao Yu², Yong Dai²,
Hongming Zhang², Zhenzhong Lan³, Dong Yu²

¹Zhejiang University, ²Tel
hehongliang@westlake.edu

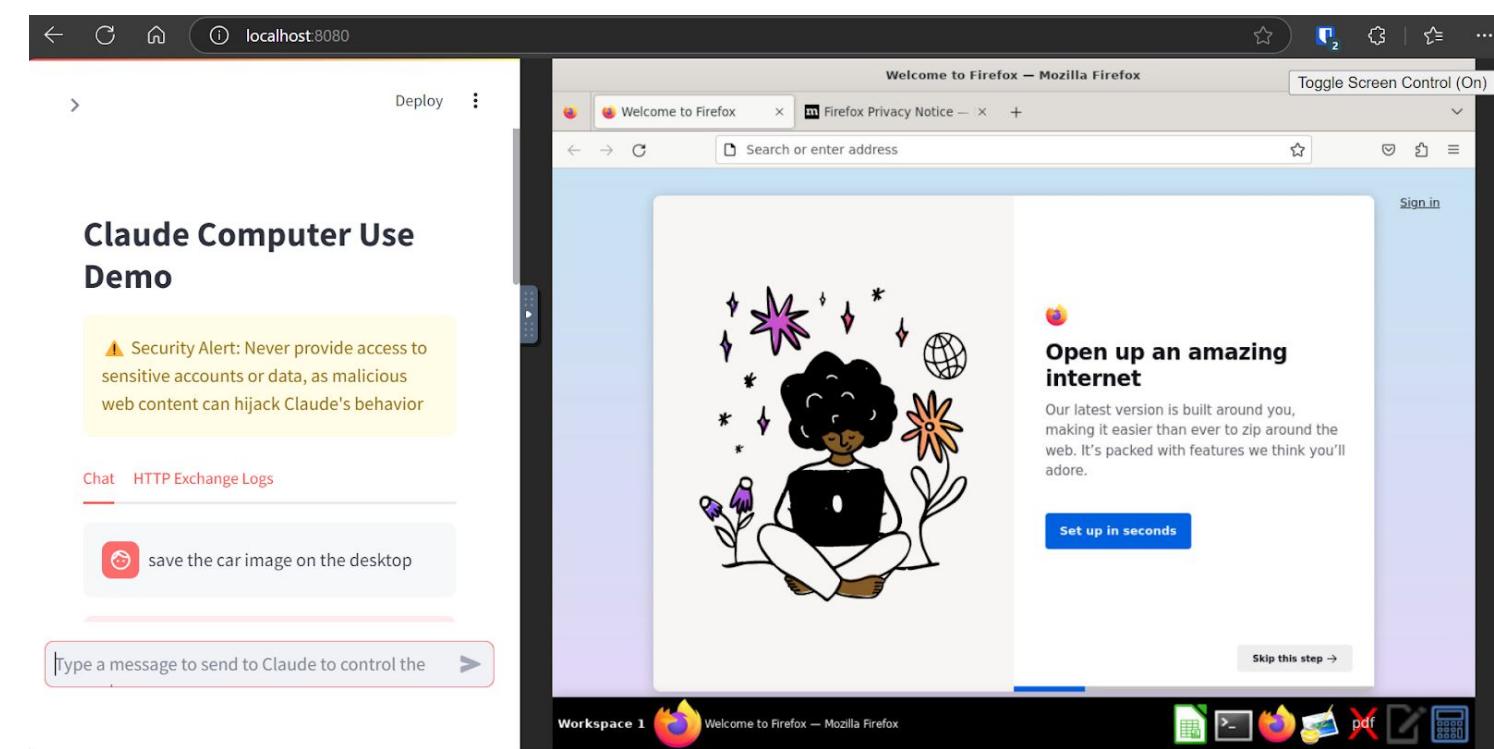
- Vision + Android OS API
- Vision + Windows OS API
- UAI + Selenium



GUI Agents - Cross Applications



Computer Use (within a Docker) by Anthropic



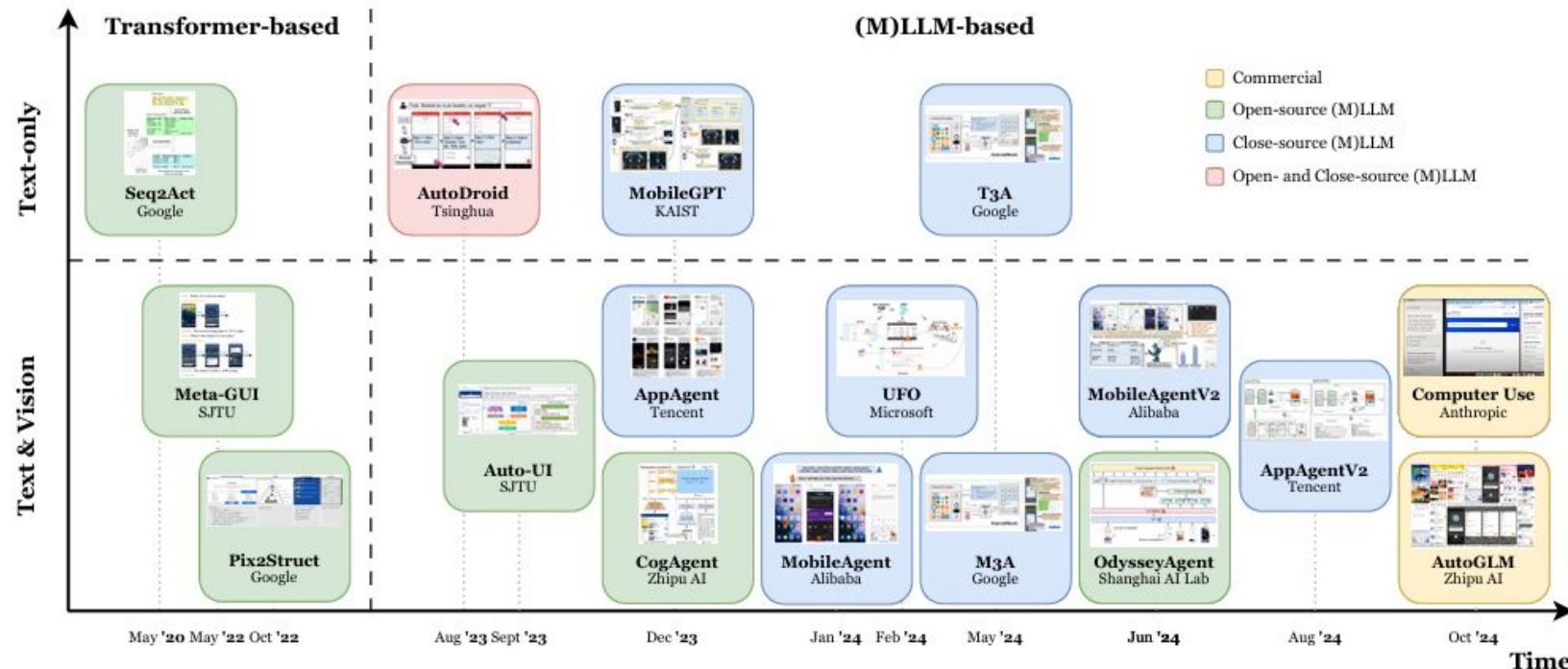
GUI Agents with Foundation Models: A Comprehensive Survey

Shuai Wang¹, Weiwen Liu^{*1}, Jingxuan Chen¹, Weinan Gan¹, Xingshan Zeng¹,
Shuai Yu¹, Xinlong Hao¹, Kun Shao¹, Yasheng Wang¹, and Ruiming Tang¹

<https://arxiv.org/abs/2411.04890>

¹Huawei Noah's Ark Lab

{wangshuai231, liuweiwen8}@huawei.com

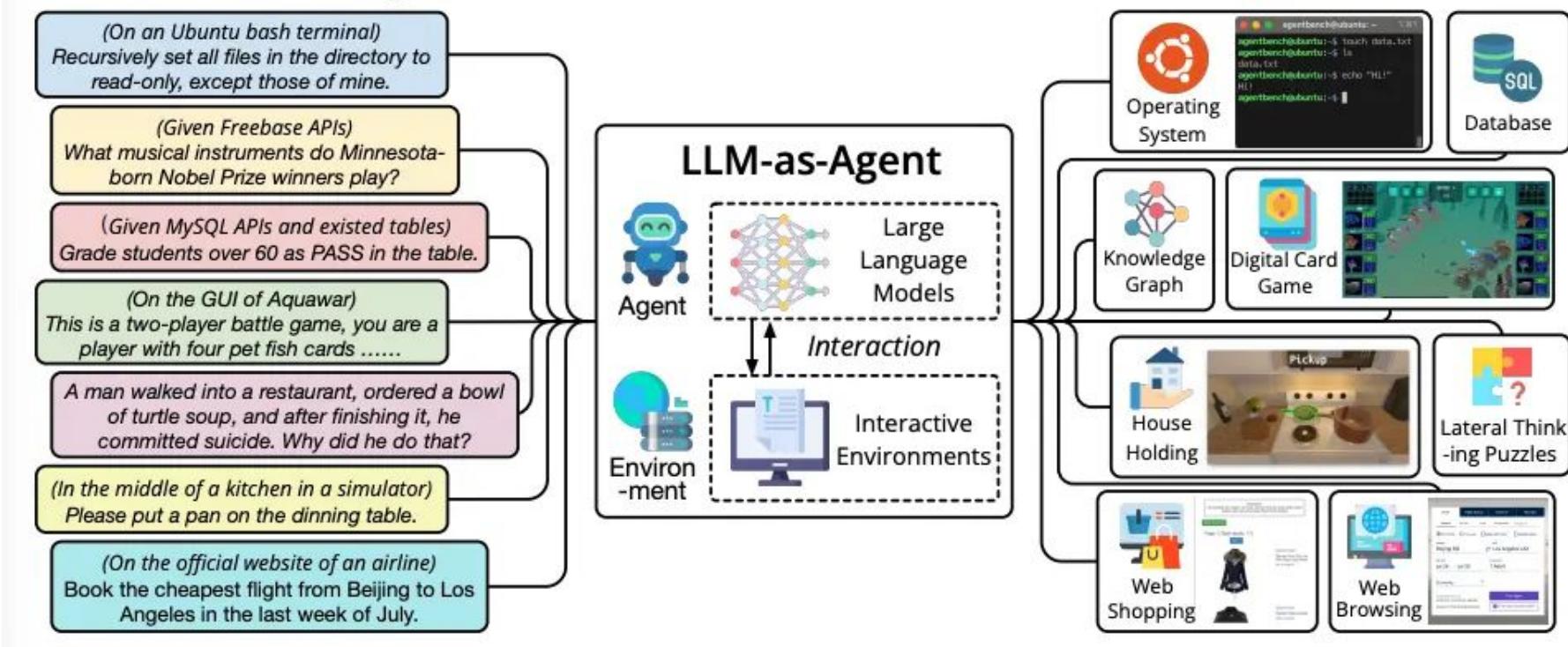


AGENTBENCH: EVALUATING LLMs AS AGENTS

Xiao Liu^{1,*}, Hao Yu^{1,*}, Hanchen Zhang¹, Yifan Xu¹, Xuanyu Lei¹, Hanyu Lai¹, Yu Gu², Hangliang Ding¹, Kaiwen Men¹, Kejuan Yang¹, Shudan Zhang¹, Xiang Deng², Aohan Zeng¹, Zhengxiao Du¹, Chenhui Zhang¹, Sheng Shen³, Tianjun Zhang³, Yu Su², Huan Sun², Minlie Huang¹, Yuxiao Dong¹, Jie Tang¹

¹Tsinghua University, ²The Ohio State University, ³UC Berkeley

<https://arxiv.org/abs/2308.03688>



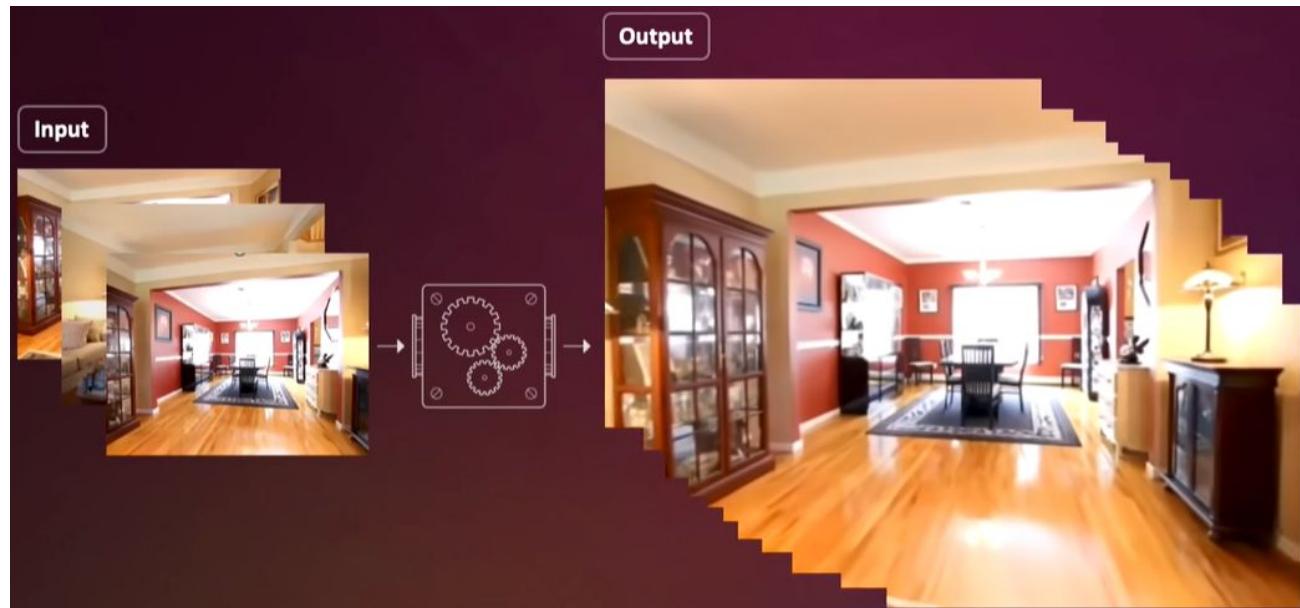
Physical World



Spatial Intelligence

- From 2D to 3D
- Seeing is for doing

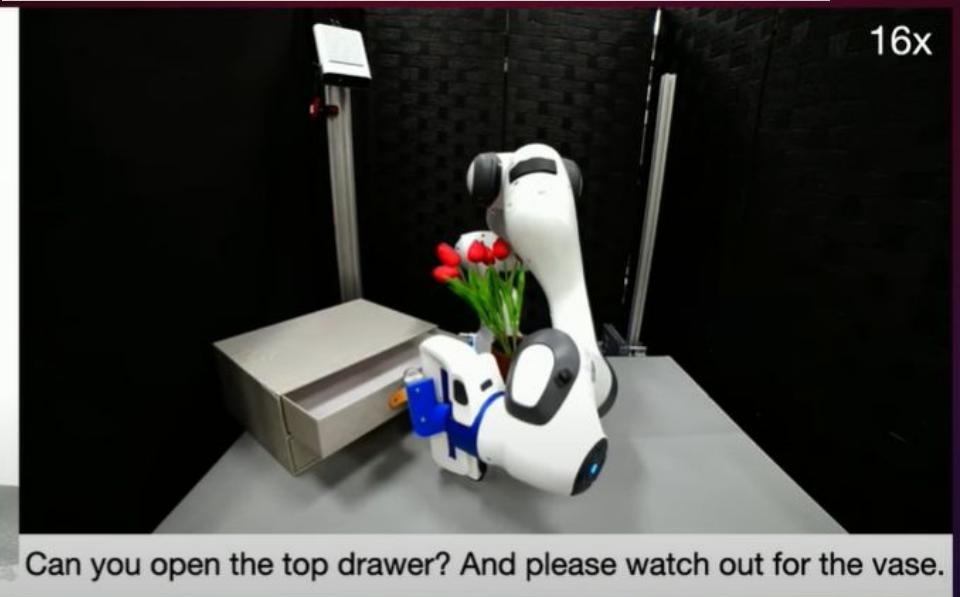
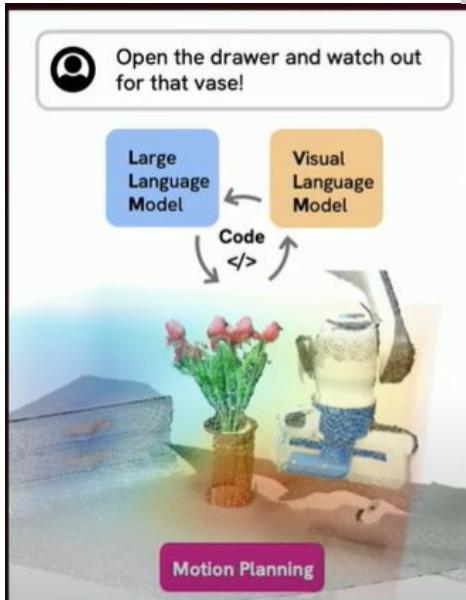
Fei Fei Li, Professor at
TED, Apr. 2024



Spatial Intelligence

- From 2D to 3D
- Seeing is for doing

Fei Fei Li, Professor | Data + AI Summit 2024



Robotic Learning

Imitation Learning
Reinforcement Learning



**CVPR24 Tutorial | Chelsea Finn:
Humanoids and Robot Generalists**



Teleoperation with ALOHA



Teleoperation with Mobile ALOHA



Text to Drive

(World Model for Automobile)

Language Augmented Driving Liquid Network

[Text to Drive | Daniela Rus |](#)

[MIT 2024](#)

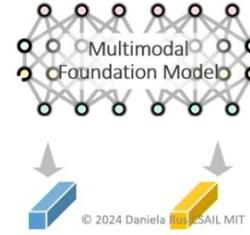


Text to Drive



Feature Vector

6/7/2024



Driving through a forest blanketed in autumn foliage, you come upon a clearing. Just ahead, a deer emerges from the trees by the roadside.

Zero-shot OOD Generalization



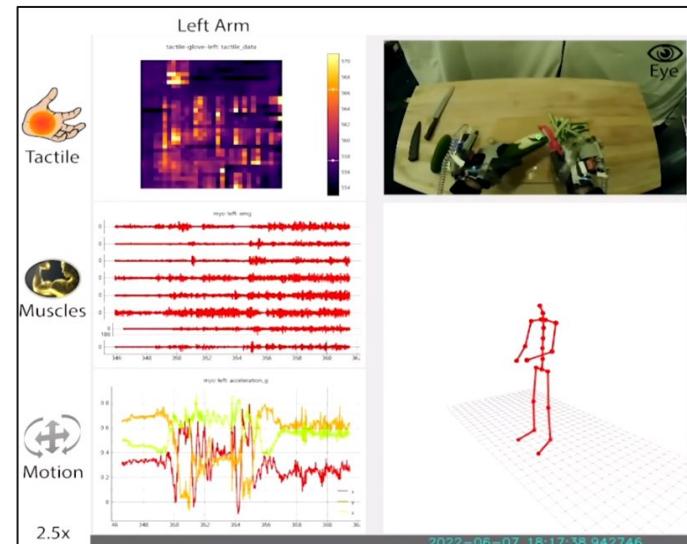
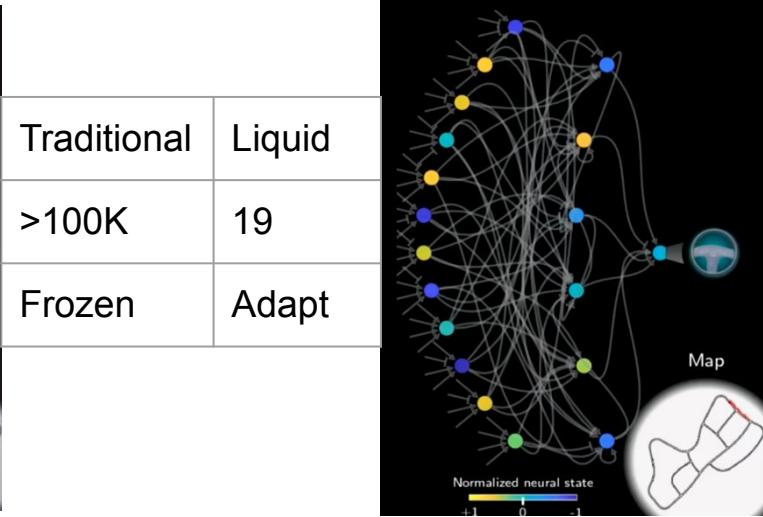
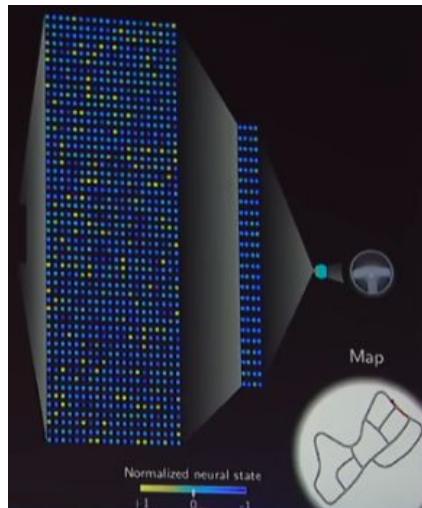
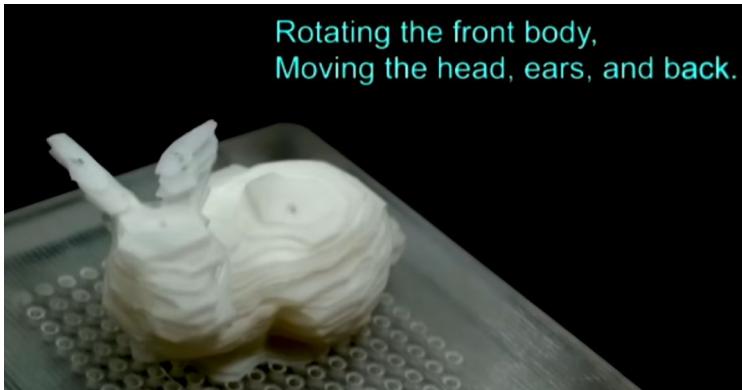
Physical Intelligence

- Think: Liquid Networks
- Design: Text/Image to Robot
- Learn from Human

**Text to Drive | Daniela Rus |
MIT 2024**



Generate design (shape, material, actuators, sensors, control program)

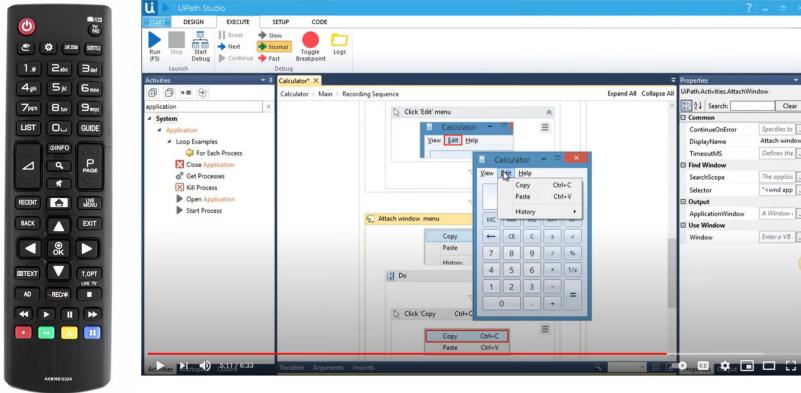




Opportunities

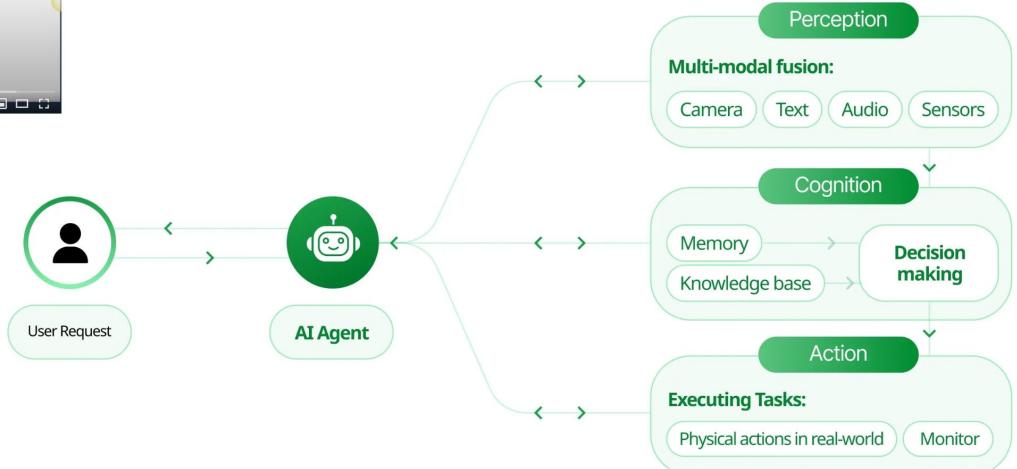
- Agentic AI
 - Civilization of RPA
 - Better HCI for all information systems
 - CUI 2025
- Sovereign and Sustainable AI
 - Data and Model Sovereignty
 - Economic Independence
 - DeepSeek R1, Grok 3
- Safety AI
 - Evaluation of AI Systems
 - Ethical and Cultural Alignment

Tools vs. Agents



With many hidden functions!

AI Agents – Brain, Perception, Tools, Agents





搜尋系統為什麼不好用？

We need better search systems

我們是不是常常在尋找 網站中的某個功能？

We need better UI and UX



會議期間同出國期間

是	會議名稱	The Thirteenth International Conference on Ubi-Media Computing (Ubi-Media 2025/ I-SPAN 2025)	
會議起日	1140119	會議迄日	1140123
*事由 <input checked="" type="checkbox"/> 同會議名稱 (不限字數)	The Thirteenth International Conference on Ubi-Media Computing (Ubi-Media 2025/ I-SPAN 2025) 擔任Ubi-Media 2025國際研討會的Keynote speaker.		
主辦或邀請單位	Kasetsart University, Thailand Tamkang University, Taiwan Santa Clara University, USA	出國期間課程安排	寒假
上傳相關證明文件	<input type="button" value="選擇檔案"/> UbiMedia 202...Hui Chang.pdf <small>(邀請函、公文、電子郵件等，每一檔案大小須為4MB以下) (檔案格式僅限: png、jpg、jpeg、gif、pdf、doc、docx，且檔案大小4MB以內。)</small>		

出國期間請假設定

假別	事實發生日/保留休假起日	起日	開始時間	迄日	結束時間
<input type="button" value="+"/> 請選擇	非 事實發生日/保留休假起日	1140119	00 : 00 該起日時間不得 小於出國起日時間	1140122	23 : 59 該迄日時間不得 大於出國迄日時間或小於出國 起日時間
請假總計(小時)		<input type="button" value="試算"/>			

Get stuck in the middle
of some process?



From Low-Code to No-Code

[https://www.nextw.com/what-is-no-code-ap](https://www.nextw.com/what-is-no-code-application)

No-Code Platforms

SME Self-Serve, Modular Interface, Immediate



Low-Code Solutions

Citizen developers, Simplified,
Cost efficient, Lower technical debt

Traditional Programming

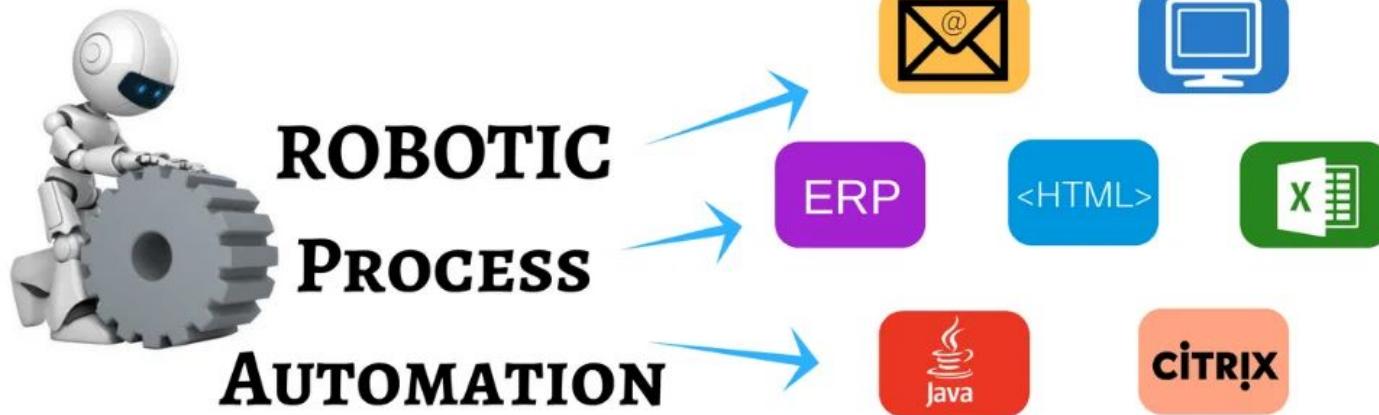
Specialists, Complex, Expensive
high technical debt

Your Turn - Discussion

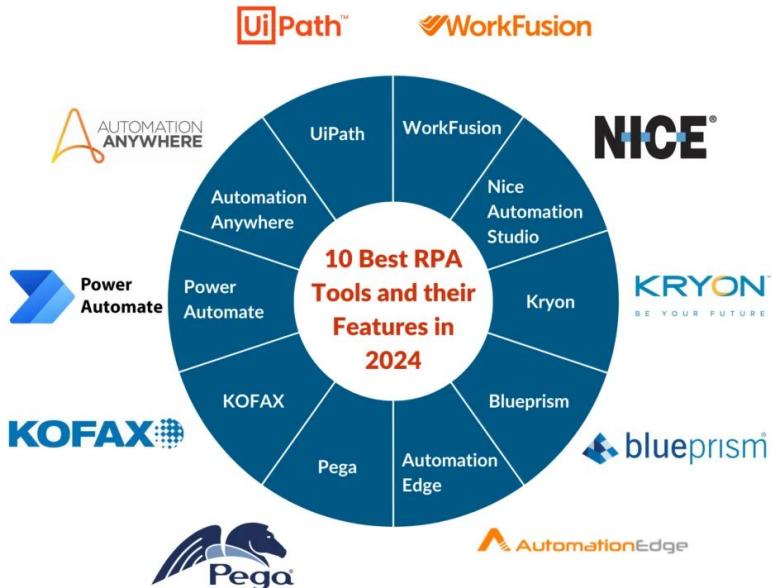
Topics:

- Difference of AI Agents and Agentic AI?
- Agentic AI in X?

Applications in RPA



RPA Tools

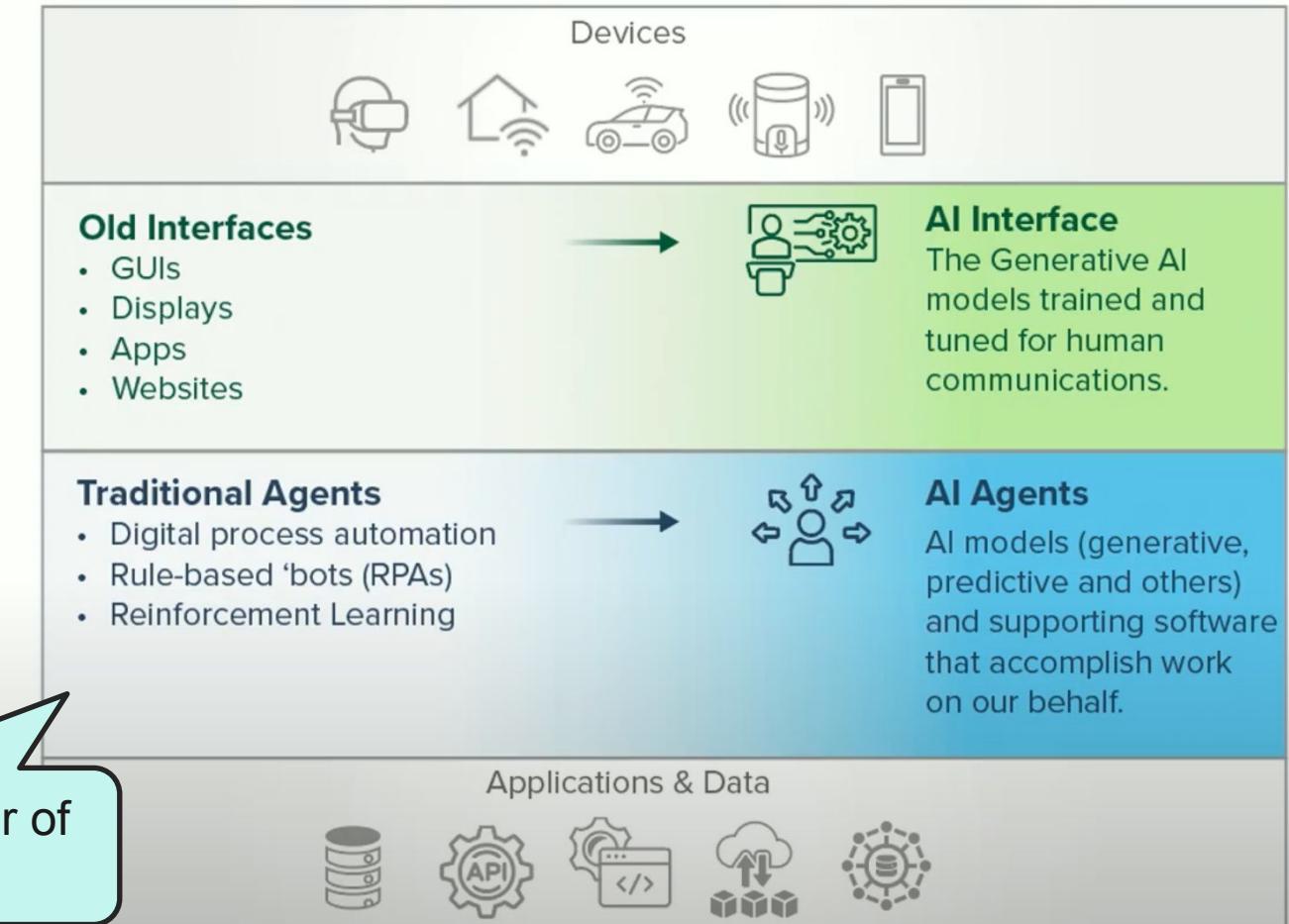


RPA creation may takes days, weeks, or months

	Automation Anywhere	ElectroNeek	UiPath	WinAutomation by Softomotive
Star Rating	★★★★★	★★★★★	★★★★★	★★★★★
Visual Editor	<div style="width: 8.7%;"></div>	<div style="width: 9.5%;"></div>	<div style="width: 8.8%;"></div>	<div style="width: 8.6%;"></div>
Meets Requirements	<div style="width: 8.7%;"></div>	<div style="width: 9.4%;"></div>	<div style="width: 8.9%;"></div>	<div style="width: 8.8%;"></div>
Ease of Use	<div style="width: 8.8%;"></div>	<div style="width: 9.6%;"></div>	<div style="width: 9.0%;"></div>	<div style="width: 8.6%;"></div>
Ease of Setup	<div style="width: 8.2%;"></div>	<div style="width: 9.5%;"></div>	<div style="width: 8.6%;"></div>	<div style="width: 9.0%;"></div>
Bot Scheduling	<div style="width: 8.9%;"></div>	<div style="width: 8.9%;"></div>	<div style="width: 8.9%;"></div>	<div style="width: 8.8%;"></div>
Attended Automation	<div style="width: 8.7%;"></div>	<div style="width: 9.8%;"></div>	<div style="width: 8.8%;"></div>	<div style="width: 9.0%;"></div>
Unattended Automation	<div style="width: 8.5%;"></div>	<div style="width: 9.6%;"></div>	<div style="width: 8.7%;"></div>	<div style="width: 8.4%;"></div>
Bot Performance Analytics	<div style="width: 8.7%;"></div>	<div style="width: 9.0%;"></div>	<div style="width: 8.7%;"></div>	<div style="width: 7.7%;"></div>



Overcome the barrier of RPA creation



Source: Forrester Report, "Change The Interface; Change The World"

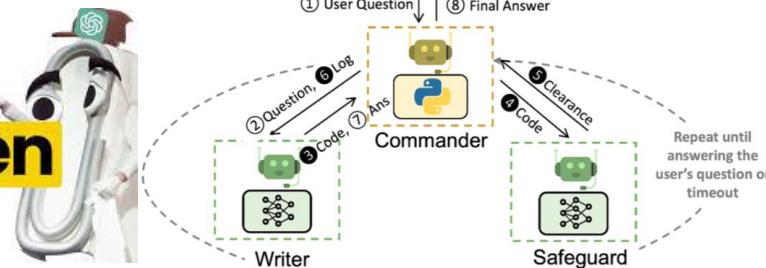
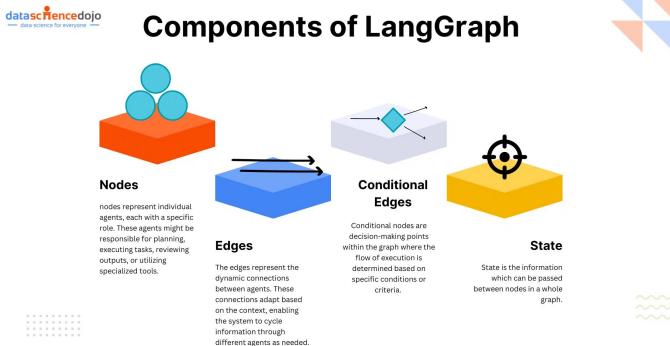
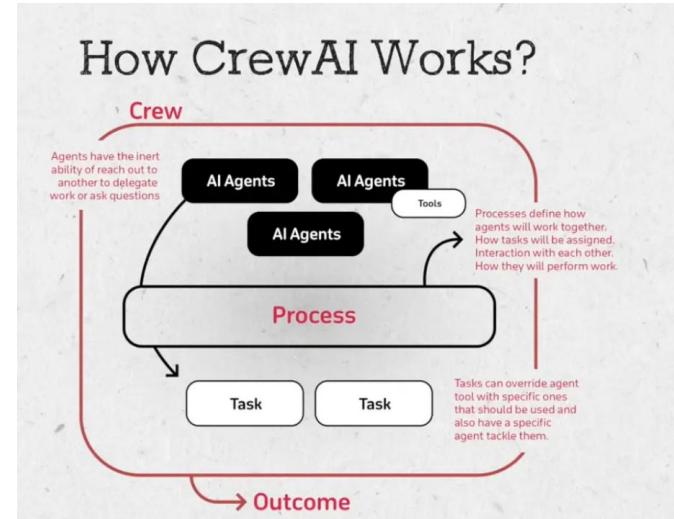
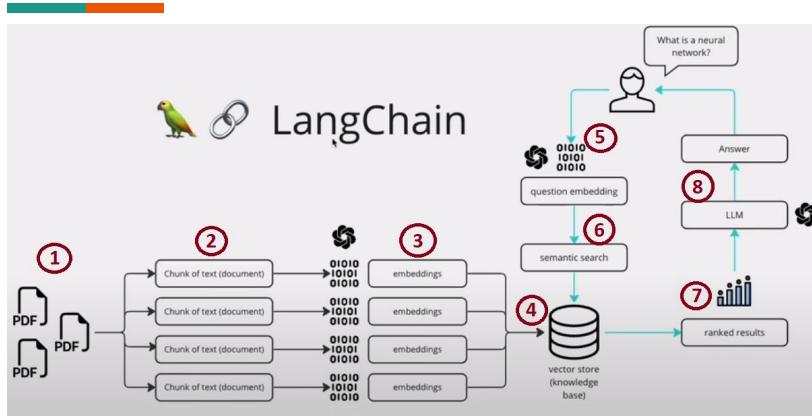
Thank you!

Thank you for your time and attention.

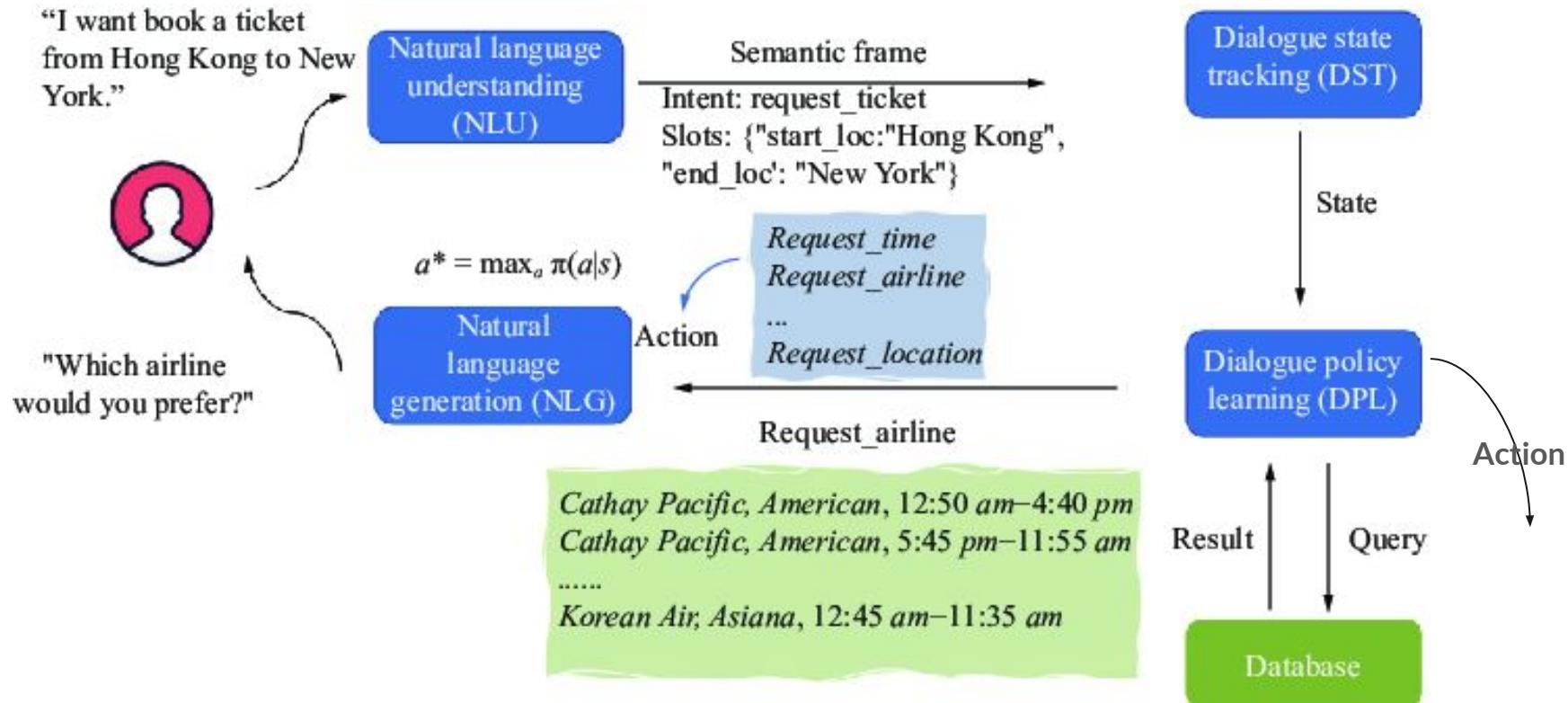
If you have any questions or would like to discuss this topic further, please feel free to email me.

chiahu@g.ncu.edu.tw

Agent Development Frameworks

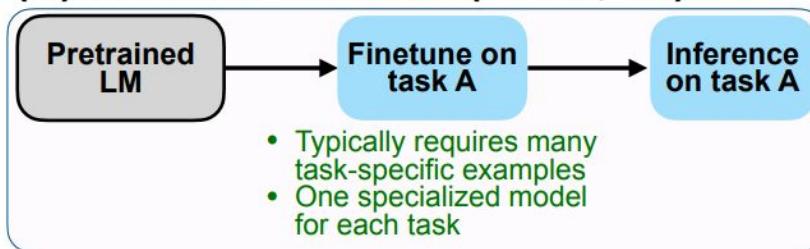


TOD: Task-Oriented Dialog Systems

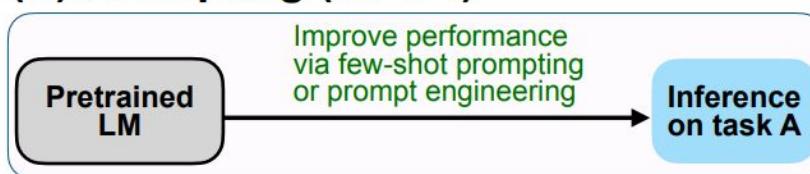


Pre-training vs. Fine-Tuning vs. In-context Learning

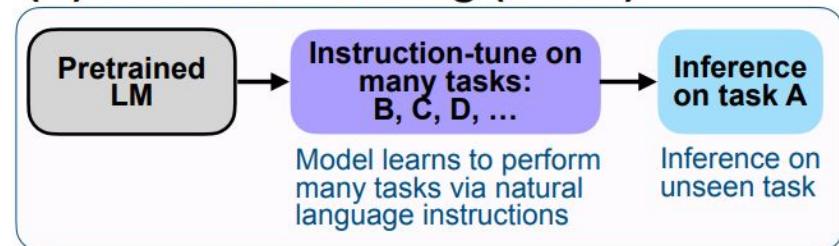
(A) Pretrain–finetune (BERT, T5)



(B) Prompting (GPT-3)



(C) Instruction tuning (FLAN)



Challenges of EduACT

1. The agent creator might have overlooked something.

The platform need to handle situations not considered by the creator.

2. Users might not always follow agents' guidances.

- Relevance score measures whether users respond well to agents
- Distance to the goal