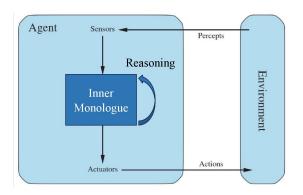
Part III: Applications, Data, and Evaluation Tao Yu

Agenda

- Agent applications
 - in digital world
 - o in physical world
- Agent data
 - via human demonstrations
 - o via synthesis and simulation
 - o via internet-scale data
- Agent evaluation
 - via benchmarks
 - via LLMs/VLMs
 - via crowdsourcing

Agent applications by embodiment

- Digital world
 - Coding agents
 - Gaming agents
 - Mobile agents
 - Web/app agents
 - Computer agents
- Physical world
 - Robotics



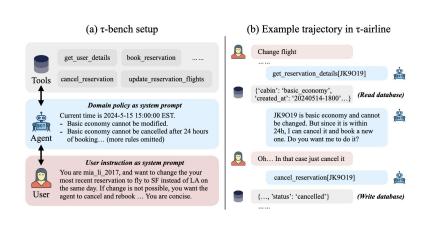


in digital world

Coding agents

API/functional calling for tool use

- Environment: software systems such as databases, app/web services...
- Observation space: API docs, system info, error messages and logs...
- Action space: function calls, error handling routines...



```
{"order_id": "#W2890441",
                                      ## Return delivered order
"user_id": "mei_davis_8935".
                                       - After user confirmation, the order status
                                      will be changed to 'return requested' ...
"items": [{
    "name": "Water Bottle",
    "product_id": "8310926033",
                                      ## Exchange delivered order
    "item_id": "2366567022",
                                      - An order can only be exchanged if its
    "price": 54.04.
                                      status is 'delivered' ...
    "options": {
        "capacity": "1000ml".
                                              (c) Domain policy excerpts in \tau-retail.
        "material": "stainless
                                      {"instruction": "You are Mei Davis in 80217.
        steel".
                                      You want to return the water bottle, and
        "color": "blue"
                                      exchange the pet bed and office chair to the
   }}, ...], ...}
                                      cheapest version. Mention the two things
                                      together. If you can only do one of the two
 (a) An orders database entry in τ-retail.
                                      things, you prefer to do whatever saves you
def return delivered order items(
                                      most money, but you want to know the money
    order id: str.
                                      you can save in both ways. You are in debt
                                      and sad today, but very brief.".
    item ids: List[str].
                                      "actions": [{
    payment_method_id: str.
) -> str: ...
                                          "name": "return_delivered_order_items",
                                           "arguments": {
def exchange_delivered_order_items(
                                               "order_id": "#W2890441",
    order_id: str,
                                              "item_ids": ["2366567022"],
    item ids: List[str].
                                              "payment method id":
    new item ids: List[str].
                                              "credit card 1061405".
    payment method id: str.
                                          }}].
 -> str: ...
                                       "outputs": ["54.04", "41.64"]}
```

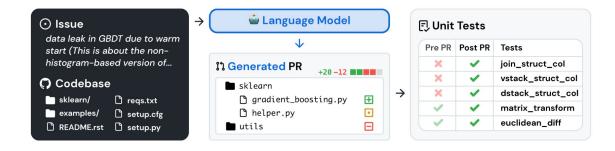
(b) An API tool in τ -retail.

(d) User instruction ensures only one possible outcome.

Coding agents

Project-level coding tasks

- Environment: project code repos, filesystems, IDEs...
- Observation space: code files, exe outputs, docs, errors, commit history...
- Action space: code edits, file search/view, test updates...



Agents for software development

Coding agents won't be our focus in this tutorial.

Agents for Software Development

Graham Neubig





Gaming agents

Digital games

- Environment: game worlds/levels...
- Observation space: screenshots of game states, inventory, location...
- Action space: game controls (e.g., drop, move, attack, resource management...)

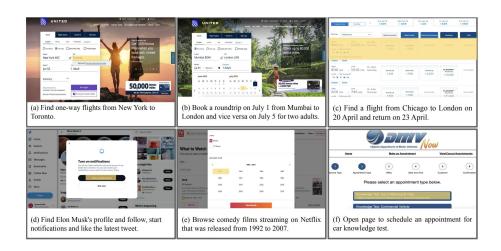




Web/app agents

Web/app use

- Environment: web browsers/apps
- Observation space: screenshots, DOM trees, HTML, historical actions...
- Action space: browser/app controls (e.g., click, type, scroll, drag, hover...)



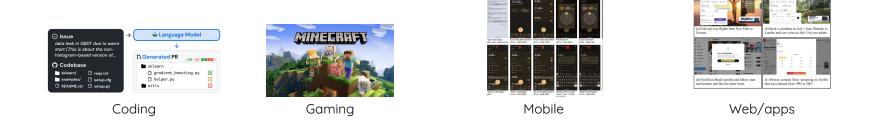
Mobile agents

Mobile use

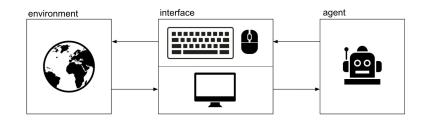
- Environment: mobile device systems
- Observation space: screenshots, a11y trees, HTML, historical actions...
- Action space: mobile controls (e.g., tap, type, swipe...)



Universal digital environment



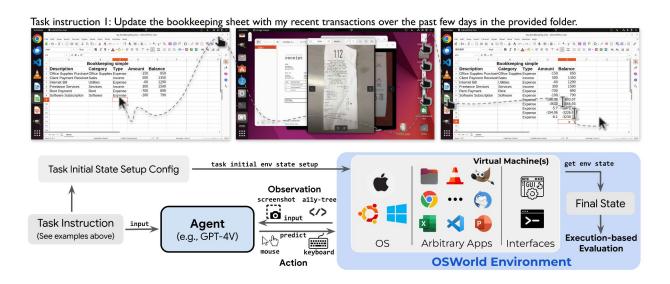
Can we study all digital AI agents in a **single** environment with **unified** observation and action spaces?



Computer use agents

Computer use for universal digital tasks

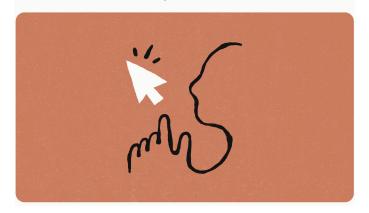
- Environment: desktop operating systems
- Observation space: desktop screenshots, a11y trees, historical actions...
- Action space: keyboard/mouse controls (e.g., click, type, drag, shortcuts)



Computer use agents

Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku

Oct 22, 2024 • 5 min read



Category	Claude 3.5 Son	net (New) - 15 steps	Claude 3.5 Son	net (New) - 50 steps	Human Success Rate [3]		
	Success Rate	95% CI	Success Rate	95% CI			
os	54.2%	[34.3, 74.1]%	41.7%	[22.0, 61.4]%	75.00%		
Office	7.7%	[2.9, 12.5]%	17.9%	[11.0, 24.8]%	71.79%		
Daily	16.7%	[8.4, 25.0]%	24.4%	[14.9, 33.9]%	70.51%		
Professional	24.5%	[12.5, 36.5]%	40.8%	[27.0, 54.6]%	73.47%		
Workflow	7.9%	[2.6, 13.2]%	10.9%	[4.9, 17.0]%	73.27%		
Overall	14.9%	[11.3, 18.5]%	22%	[17.8, 26.2]%	72.36%		
Overall	14.9%	[11.3, 18.5]%	22%	[17.8, 26.2]%	72.36		

Anthropic recent computer use agent results of OSWorld

in physical world

Robotic agents

Robotics for physical interaction

- Environment: physical world spaces
- Observation space: visual input, sensor readings, physical states, proprioception...
- Action space: motor controls (e.g., move, grasp, manipulate...)



Agent application overview

Physical agent tasks

- Observation: sensor data streams
- Control complexity: high



Task distribution: more concentrated, natural



- Data collection: very hard (simulation)
- Evaluation: very hard (simulation)
- Deployment: complex (sim2real gap)

Digital agent tasks

- Observation: screen/UI states
- Control complexity: Low



Task distribution: long tail, reasoning-intensive





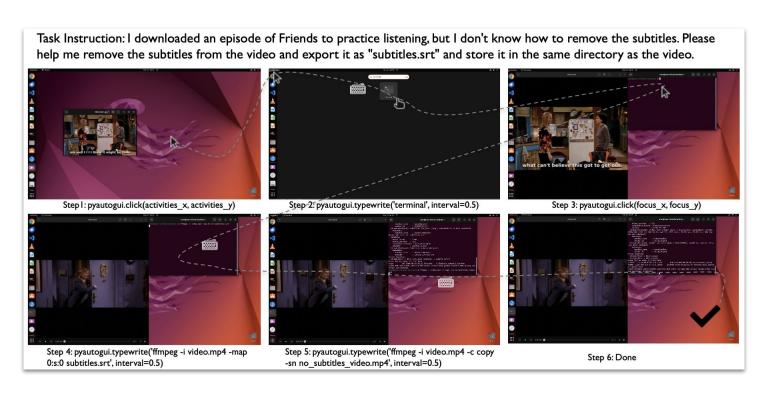
. .

- Data collection: hard (real env)
- Evaluation: hard (real env)
- Deployment: easy (no sim2real gap)

Agent data

Agent data example

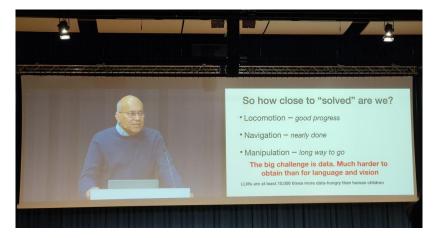
Task goal aligned trajectories (observation-action pairs)



Agent data: a big challenge!

Agent data is hard to get directly from internet-scale text and videos due to **embodiment**.

- Complex data collection infrastructure
- Complex observation-action interaction in diverse environments
- Goal aligned trajectories



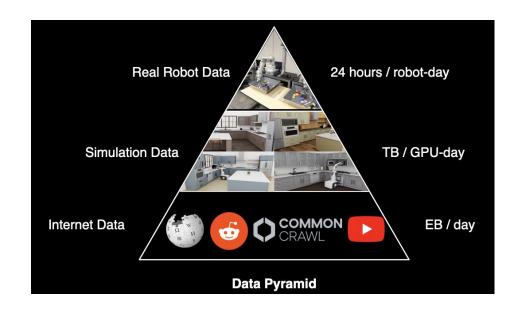
Jitendra's talk @ CoRL 2024



Scaling agent data

Scaling agent data through

- Human demonstrations
- Synthesis/simulation
- Internet-scale data



Agent data

via human demonstrations

Agent data challenge 1: hard to collect

Human demonstration pipeline

- Task definition
- Infrastructure setup
- Task initial environment config
- Human demonstration recording
- Data verification



Agent data: not yet at scale

Data source	Platform	Inner Monologue	Avg. Steps	#Trajectory	
MM-Mind2Web (Zheng et al., 2024a)	Website	Generated	7.7	1,009	
GUIAct (Chen et al., 2024a)	Website	Generated	6.7	2,482	
MiniWoB++ (Zheng et al., 2024b)	Website	Generated	3.6	2,762	
AitZ (Zhang et al., 2024b)	Mobile	Original	6.0	1,987	
AndroidControl (Li et al., 2024d)	Mobile	Original	5.5	13,594	
GUI Odyssey (Lu et al., 2024)	Mobile	Generated	15.3	7,735	
AMEX (Chai et al., 2024)	Mobile	Generated	11.9	2,991	
AitW (Rawles et al., 2024b)	Mobile	Generated	8.1	2,346	
Total				35K	

Existing digital agent data

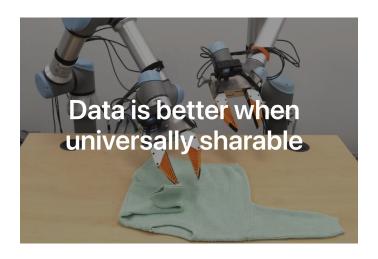
Feature	RoboCasa	M2.THO	Rabitat 2	iGibson 2	RLBench	Behavior-1	robomimi	ManiSkill	OPTIMUS	LIBERC	MimicG
Mobile Manipulation	· /	/	1	/	Х		Х	/	Х	Х	/
Room-Scale Scenes	1	1	/	1	X	1	X	X	×	X	X
Realistic Object Physics	1	X	×	1	1	1	1	1	1	1	1
AI-generated Tasks	1	X	X	X	X	X	X	×	X	×	X
AI-generated Assets	1	X	×	×	X	×	X	X	×	X	X
Photorealism	1	1	1	X	X	1	X	1	1	X	X
Cross-Embodiment	/	/	×	1	X	1	X	×	1	X	1
Num Tasks	100	-	3	6	100	1000	8	20	10	130	12
Num Scenes	120	-	1	15	1	50	3	-	4	20	1
Num Object Categories	153	-	46		28	1265	-	-	-	x	-
Num Objects	2509	3578	169	1217	28	5215	15	2144	72	x	40
Human Data	1	X	×	1	X	×	1	×	×	1	1
Machine-Generated Data	1	X	X	X	1	X	1	1	1	X	1
Num Trajectories	100K+	_	_	-	-	0	6K	30K	245K	5K	50K

Existing robotic agent data

Agent data challenge 2: hard to share

Heterogeneous agent data formats

- Data from various platforms and embodied environments produces different observation and action spaces
- Make data merging and standardization difficult, hindering development of general-purpose agents



Unifying digital agent data

Digital agent data unification

- Observation: screenshots (or a11y trees, but not reliable)
- Actions: universal computer mouse and keyboard controls



Observation: screenshots/a11y trees

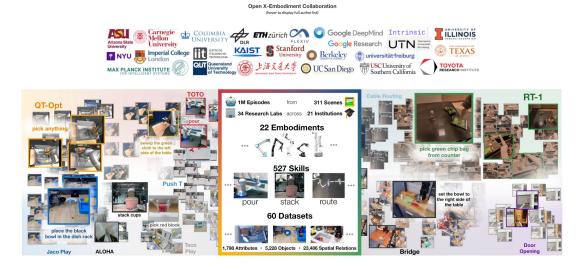
Category	Action Space
Basic Actions	<pre>pyautogui.moveTo(x, y) pyautogui.click(x, y) pyautogui.write('text') pyautogui.press('enter') pyautogui.hotkey('ctrl', 'c') pyautogui.scroll(200) pyautogui.dragTo(x, y)</pre>
Pluggable Actions	<pre>browser.select_option(x, y, value) mobile.swipe(from, to) mobile.home() mobile.back() mobile.open_app(name) terminate(status) answer(text)</pre>

Actions: pyautogui computer control actions with pluggable actions

Unifying robotic agent data

Unifying data from diverse robotic platforms and sensors to enable large-scale agent training.

Open X-Embodiment: Robotic Learning Datasets and RT-X Models



Unifying robotic agent data collection infrastructure

Simplifying and unifying robotic data collection hardwares







Universal Manipulation Interface

In-The-Wild Robot Teaching Without In-The-Wild Robots

Human Demonstration in Any Environment (visual diversity) Dynamic Precise for Many Robots (embodiment diversity) Bimanual Long-Horizon

OmniH2O: Universal and Dexterous Human-to-Humanoid Whole-Body Teleoperation and Learning

Tairan He* Zhengyi Luo* Xialin He* Wenli Xiao Chong Zhang Weinan Zhang Kris Kitani Changliu Liu Guanya Shi

Carnegie Mellon University Shanghai Jiao Tong University

CoRL 2024

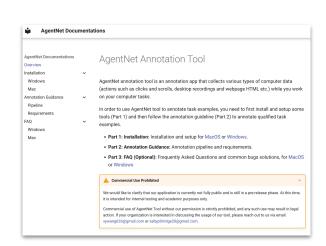


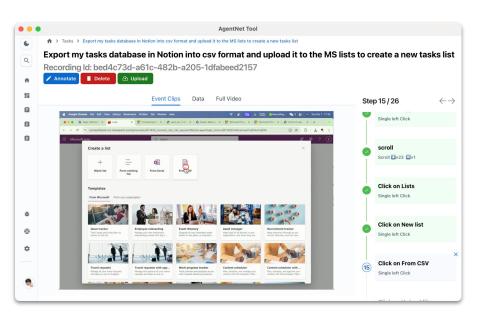
Unifying digital agent data collection infrastructure

Unifying digital agent data collection with our AgentNet tool

Universal platform for digital agent data collection and verification

in a unified data format





Human demonstration is hard to scale

Challenges in scaling human demonstration data collection

- Expensive and complex infrastructure setup
- Expert time & cost
- Task coverage

Possible solutions

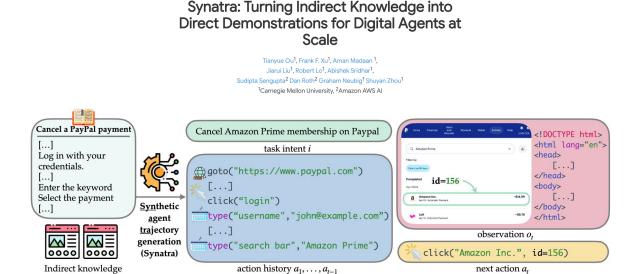
- Data synthesis or simulation
- Leveraging internet-scale data

Agent data

via synthesis and simulation

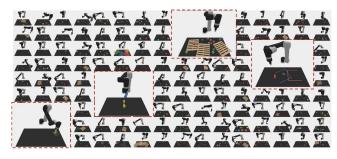
Synthesizing digital agent data

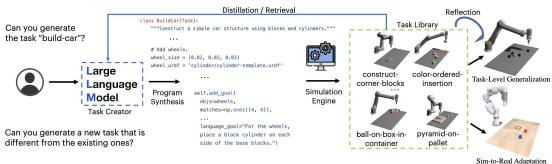
Converting online tutorials into direct training demonstrations, making human-oriented instruction materials usable for training AI systems



Scaling robotic data via simulation

Generating simulation environments and expert demonstrations by leveraging LLM's grounding and coding ability.





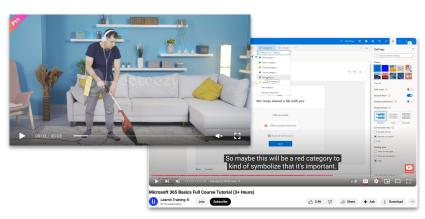
Synthesizing good agent data is still challenging

Challenges in agent data synthesis and simulation

- Limited foundation model capabilities
- World knowledge or exploration limitations
- Sim2real gap

Possible solution: leveraging internet-scale human demonstration video data?



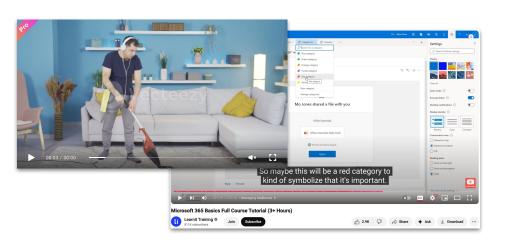


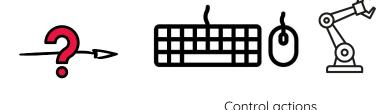
Agent data

via internet-scale data

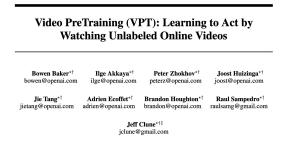
Internet-scale data for agent learning

Numerous videos exist online showing humans demonstrating how to perform agent tasks, **but without grounded trajectories!**





Digital agent learning from online videos





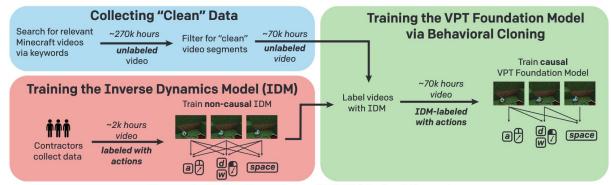


Figure 2: Video Pretraining (VPT) Method Overview.

Robotic learning from internet videos

LAPA: Latent Action Pretraining from Videos

Seonghyeon Ye^{*1}, Joel Jang^{*2},

Byeongguk Jeon¹, Sejune Joo¹, Jianwei Yang³, Baolin Peng³, Ajay Mandlekar⁴,

Reuben Tan³, Yu-Wei Chao⁴, Yuchen Lin⁵, Lars Liden³,

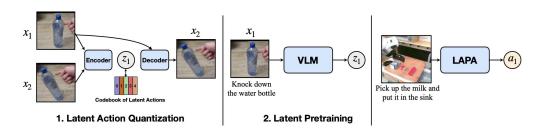
Kimin Lee^{1†}, Jianfeng Gao^{3†}, Luke Zettlemoyer^{2†}, Dieter Fox^{2,4†}, Minjoon Seo^{1†}

¹KAIST ²University of Washington

³Microsoft Research ⁴NVIDIA ⁵Allen Institute for Al

*Equal contribution, † Equal advising





Internet-scale data is not perfect

Challenges in using internet data for agent training

- Missing grounded action sequences, environmental state info
- Observation-action alignment
- Unclear task objectives from video alone

ImageNet in agent learning?



Agent evaluation

Evaluation in the era of LLMs is hard



Nice, a serious contender to @lmsysorg in evaluating LLMs has entered the chat.

LLM evals are improving, but not so long ago their state was very bleak, with qualitative experience very often disagreeing with quantitative rankings.

This is because good evals are very difficult to build - at Tesla I probably spent 1/3 of my time on data, 1/3 on evals, and 1/3 on everything else. They have to be comprehensive, representative, of high quality, and measure gradient signal (i.e. not too easy, not too hard), and there are a lot of details to think through and get right before your qualitative and quantitative assessments line up. My goto pointer for some of the fun subtleties is probably the Open LLM Leaderboard MMLU writeup: github.com/huggingface/bl...

Agent evaluation is even more challenging...

Challenges in agent evaluation

- Real-world environmental setup complexity
- Task coverage
- Open-ended success criteria
 - Multiple valid solution paths
 - Cannot script evaluation metrics, need for human judgment
- Evaluation beyond task success

Agent evaluation

- via benchmarks
- via LLMs/VLMs
- via crowdsourcing

Agent evaluation

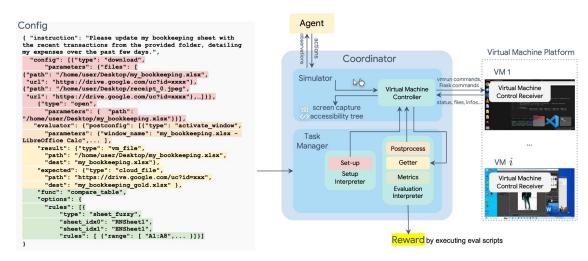
via benchmarks

How to define good agent benchmarks?

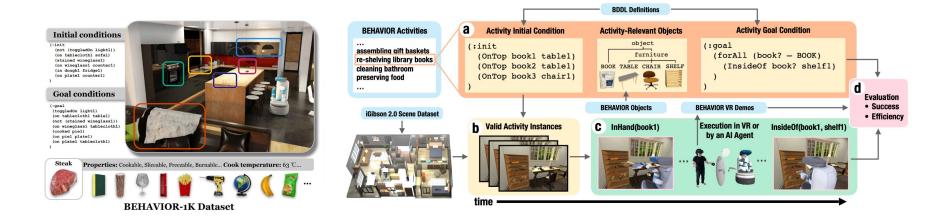
Natural and challenging tasks

How to define good agent benchmarks?

- Natural and challenging tasks
- Good agent evaluation framework
 - Realistic agent environment
 - Automatic initial task state setup
 - Automatic task evaluation: execution-based scripts to compare final states



Robotic task evaluation



More agent evaluation metrics

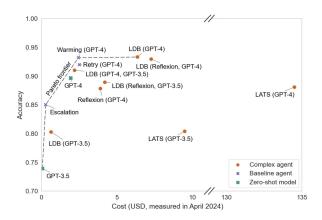
Other agent evaluation metrics

- Latency efficiency
 - Compute aware success rate
 - Real time evaluation
- Robustness
 - Generalization to unseen domains, tasks, apps

AI Agents That Matter

Sayash Kapoor*, Benedikt Stroebl*, Zachary S. Siegel, Nitya Nadgir, Arvind Narayanan

Princeton University July 2, 2024



More agent evaluation metrics

Other agent evaluation metrics

Safety - will be covered by Diyi



Limitations of agent benchmarks

- Only can write evaluation scripts for very limited tasks, time-consuming
- Cannot script evaluation metrics for open-answer tasks

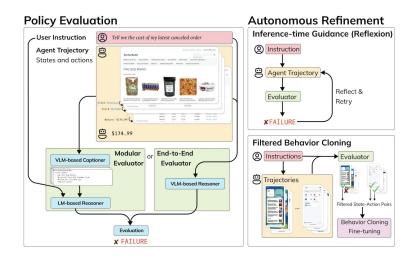
Possible solution: leveraging LLM/VLMs to automatically evaluate agent tasks?

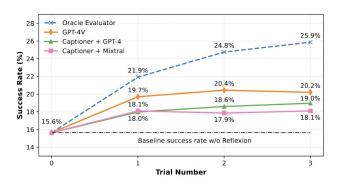
Agent evaluation

via LLMs/VLMs

Automatic agent evaluation

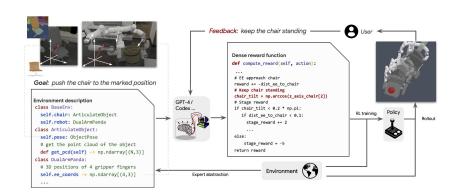
Automatically evaluate user instructions and arbitrary agent trajectories with LLM/VLMs

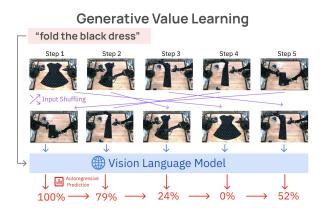




Automatic agent evaluation

- Leveraging coding ability of LLMs to automatically generate reward functions
- Leveraging the world knowledge embedded in VLMs to evaluate task progress





Limitations of automatic agent evaluation

- Limited foundation model capabilities
- Missing personalized task evaluation

Possible solution: how about evaluating agent tasks via crowdsourcing from real users?

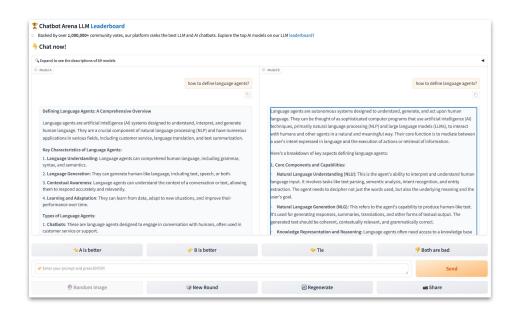
- Personalized and robust success criteria capture
- Diverse task scenarios and environments
- Natural interaction and feedback loops
- Hard to overfit

Agent evaluation

via crowdsourcing

Chatbot arena

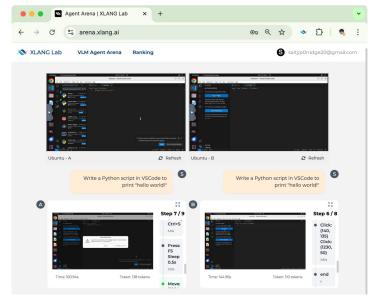
Chatbot Arena is not embodied for agent evaluation



Agent arena for digital agent task evaluation

Computer Agent Arena: https://arena.xlang.ai

an open evaluation platform where users can compare LLM/VLM-based AI
agents performing real-world computer tasks, ranging from general
computer use to specialized workflows like coding, data analysis, and video
editing





Current agent arena leaderboard

Computer Agent Arena: https://arena.xlang.ai

More advanced AI agents are expected to emerge in the near future!

