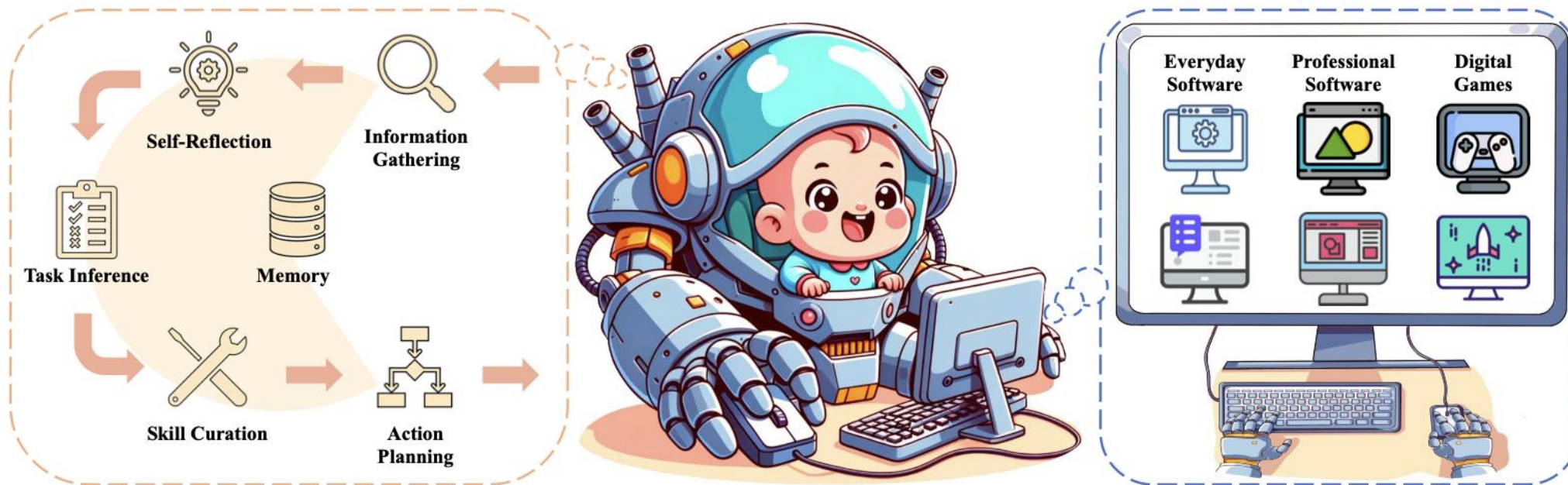


# CRADLE: Empowering Foundation Agents Towards General Computer Control



Weihao Tan<sup>3\*</sup>, Wentao Zhang<sup>3\*</sup>, Xinrun Xu<sup>5\*</sup>, Haochong Xia<sup>3†</sup>, Ziluo Ding<sup>2†</sup>, Boyu Li<sup>3†</sup>, Bohan Zhou<sup>4†</sup>,  
Junpeng Yue<sup>4†</sup>, Jiechuan Jiang<sup>4†</sup>, Yewen Li<sup>3†</sup>, Ruyi An<sup>3†</sup>, Molei Qin<sup>3†</sup>, Chuqiao Zong<sup>3†</sup>, Longtao Zheng<sup>3†</sup>,  
YuJie Wu<sup>1†</sup>, Xiaoqiang Chai<sup>1†</sup>, Yifei Bi<sup>2</sup>, Tianbao Xie<sup>6</sup>, Pengjie Gu<sup>3</sup>, Xiyun Li<sup>2</sup>, Ceyao Zhang<sup>7</sup>, Long Tian<sup>1</sup>,  
Chaojie Wang<sup>1</sup>, Xinrun Wang<sup>3‡</sup>, Börje F. Karlsson<sup>2‡</sup>, Bo An<sup>3,1§</sup>, Shuicheng Yan<sup>1§</sup>, Zongqing Lu<sup>4,2§</sup>

<sup>1</sup> Skywork AI <sup>2</sup> Beijing Academy of Artificial Intelligence <sup>3</sup> Nanyang Technological University, Singapore

<sup>4</sup> Peking University <sup>5</sup> Institute of Software, Chinese Academy of Sciences <sup>6</sup> The University of Hong Kong

<sup>7</sup> The Chinese University of Hong Kong, Shenzhen

\* Equal contribution † Core contribution ‡ Equal advising § Corresponding authors

# General Computer Control (GCC)



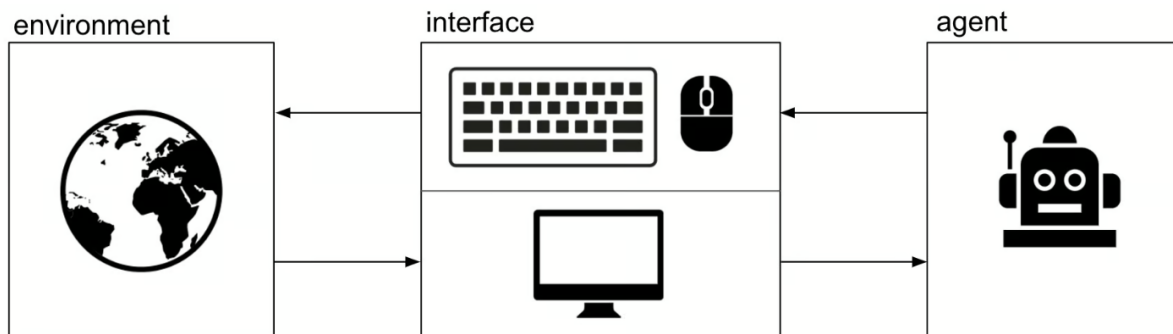
## Advantages

NO built-in APIs needed

Generalization to (unseen) tasks

Data collection with the same granularity for self-improvement

Building foundation agents that can master **ANY** computer task via the universal human-style interface by receiving input from **screens** and **audio** and outputting **keyboard** and **mouse** actions.



## Challenges

Good Alignment across multimodal inputs

Precise control of keyboard & mouse

Long-term memorizing & reasoning

Efficient exploration and self-improving

# Cradle Framework

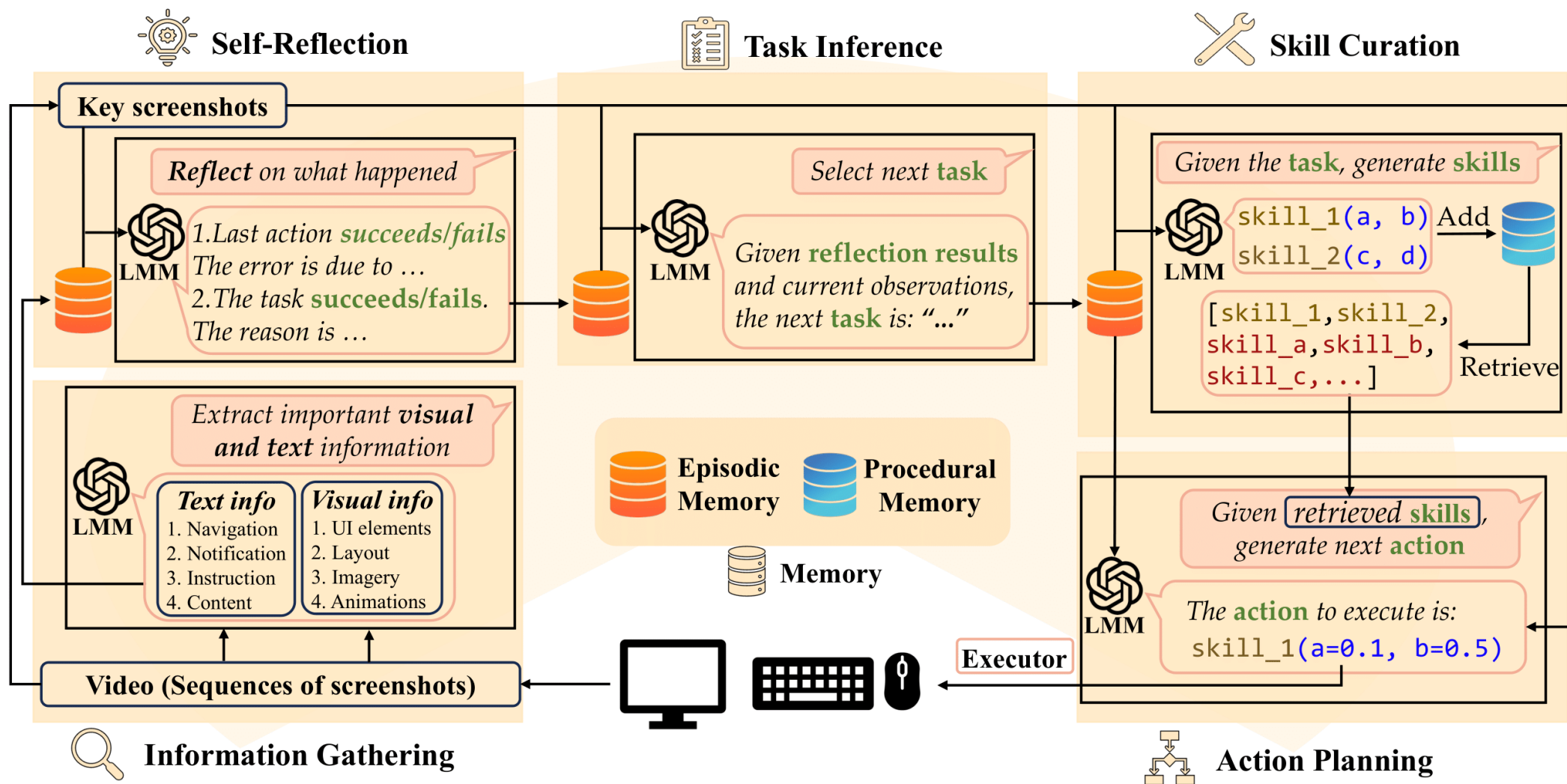
➤ Cradle: A Novel Agent Framework under GCC Setting



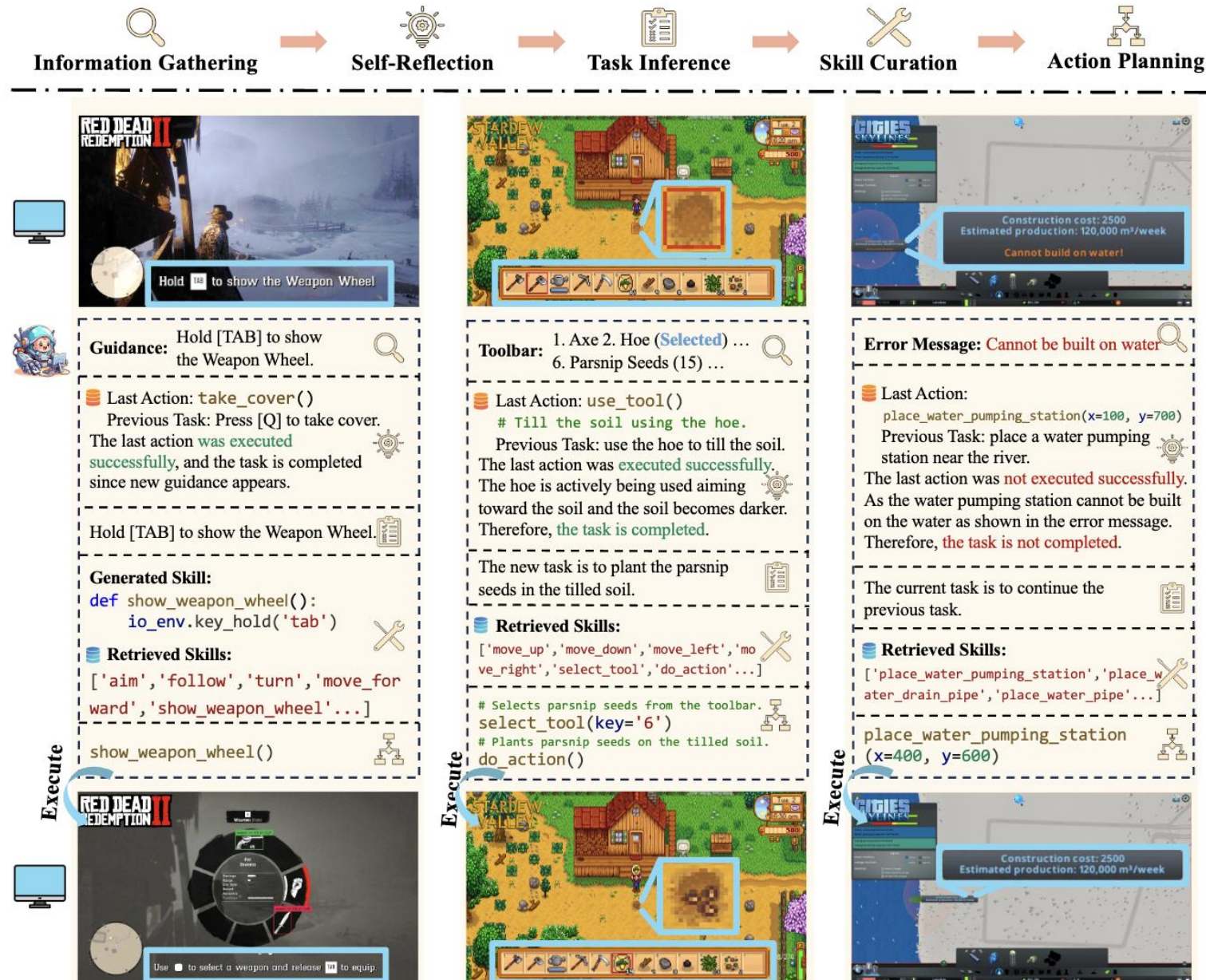


# Cradle Framework

## ➤ Six key modules



# Complete Workflow

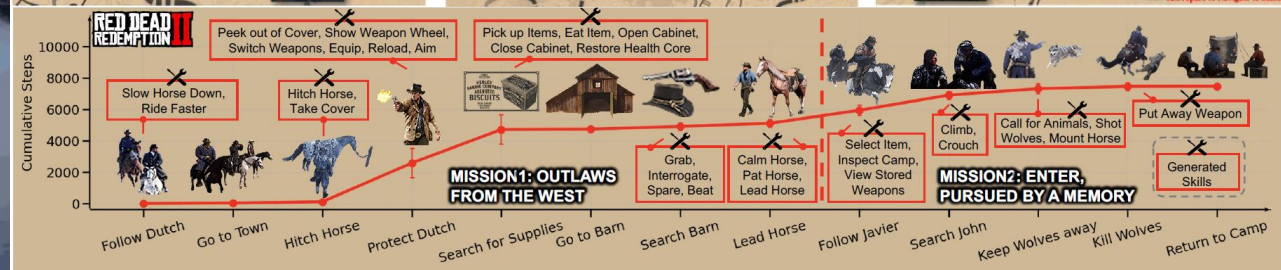






## Cradle in RDR2

- 60 minutes main storyline missions
- Long-horizon open-ended tasks







## Cradle in Other Games

- **Stardew Valley**
  - ❑ Farm clearup, cultivation and shopping
- **Dealer's Life**
  - ❑ Run a pawn shop for a week and maximize profits
- **Cities: Skylines**
  - ❑ Build a reasonable city and maximize the population with initial budgets

Stardew Valley		
Task	CRADLE	Human
Farm Clearup	14.8 ± 5.0	35.2 ± 14.5
Grids Cultivation	4/5	5/5
Shopping	1/5	5/5

Dealer's Life 2		
Metrics	CRADLE	Human
Avg. Hagglng Count	1.95 ± 0.43	1.63 ± 0.53
Turnover Rate (%)	93.6 ± 6.9	68.4 ± 22.2
Item Profit Rate (%)	37.8 ± 19.1	21.1 ± 13.6
Total Profit Rate (%)	39.6 ± 27.3	17.3 ± 15.1

Cities: Skylines		
Metrics	CRADLE	Human
Smooth Closed-loop Road	4/5	5/5
Sufficient Water Supply	1/5	3/5
Sufficient Power Supply	5/5	5/5
High Zoning Coverage	4/5	4/5
Population	450 ± 224	415 ± 416

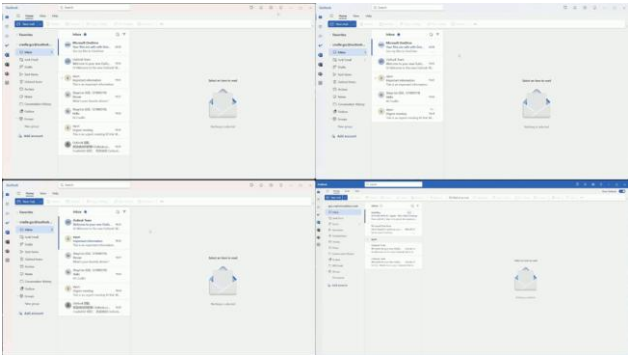


# Cradle in Software Applications

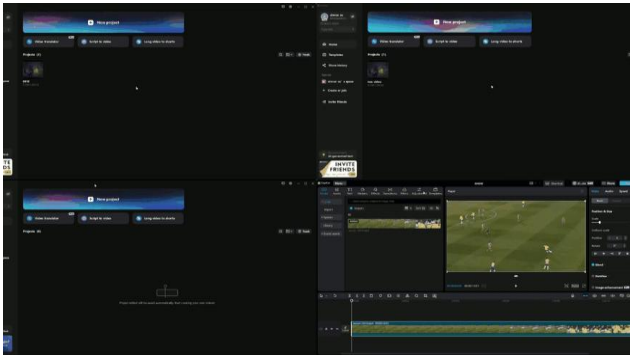
Chrome



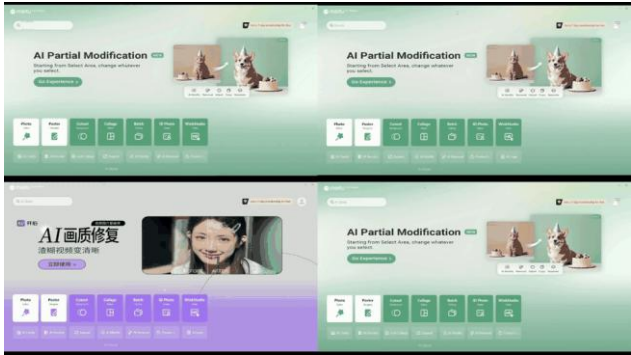
Outlook



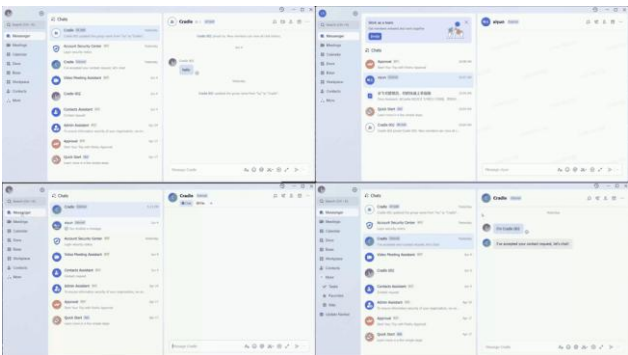
CapCut



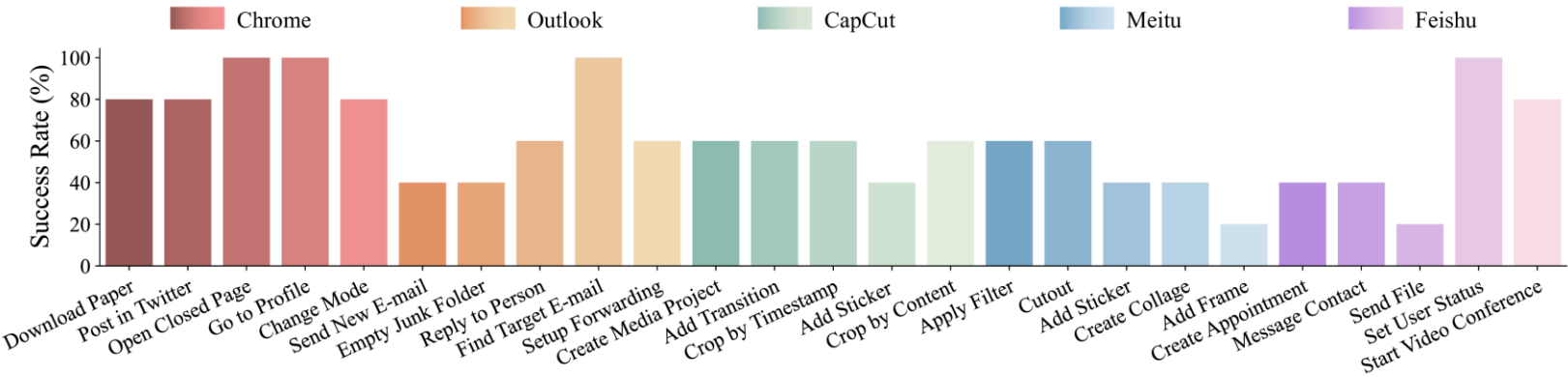
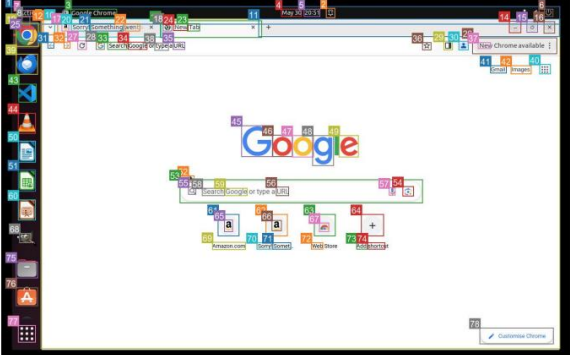
Meitu Xiuxiu



Feishu



SAM2SOM Labels



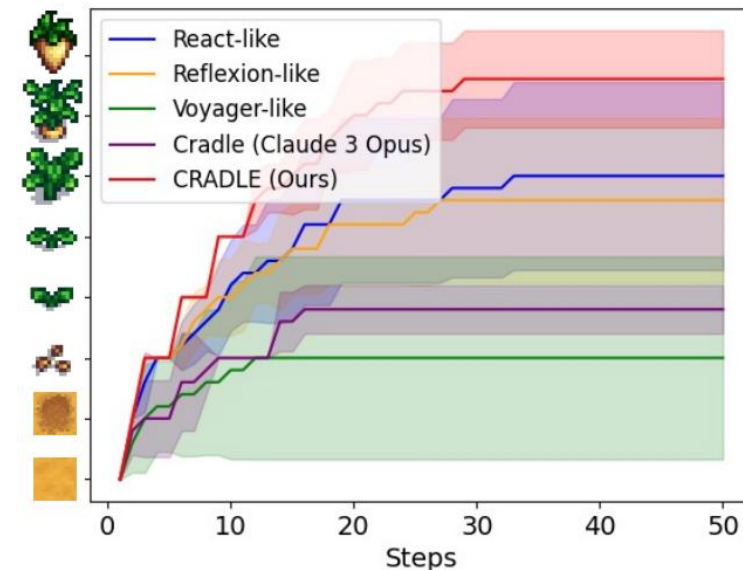
Cradle in OSWorld

Method	Office (117)	OS (24)	Daily (78)	Workfl-ow(101)	Professi-onal (49)	All (369)
GPT-4o	3.58	8.33	6.07	5.58	4.08	5.03
GPT-4o+SoM	3.58	20.83	3.99	3.60	2.04	4.59
CRADLE	3.58	16.67	6.55	5.48	20.41	7.81



# Ablation Studies & Baselines Comparison

Method	Follow Dutch	Follow Micah	Hitch Horse	Protect Dutch	Search for Supplies	Cultivation
React [72]-like (GPT-4o)	$15 \pm 2$ (5/5)	$74 \pm 0$ (1/5)	N/A	N/A	N/A	N/A
Reflexion [53]-like (GPT-4o)	$19 \pm 4$ (5/5)	$58 \pm 14$ (2/5)	N/A	N/A	N/A	N/A
Voyager [57]-like (GPT-4o)	$32 \pm 12$ (3/5)	N/A	N/A	N/A	N/A	N/A
CRADLE (Claude 3 Opus)	$30 \pm 7$ (5/5)	$52 \pm 17$ (4/5)	N/A	N/A	N/A	N/A
<b>CRADLE (GPT-4o) (Ours)</b>	<b><math>13 \pm 3</math> (5/5)</b>	<b><math>33 \pm 3</math> (5/5)</b>	<b><math>26 \pm 5</math> (4/5)</b>	<b><math>461 \pm 0</math> (1/5)</b>	<b><math>134 \pm 0</math> (1/5)</b>	<b><math>24 \pm 4</math> (4/5)</b>



Subtask	w/o Information Gathering	w/o Self-Reflection	w/o Task Inference	w/o Episodic Memory	<b>CRADLE</b>
Follow Micah	0%	0%	40%	80%	<b>100%</b>
Hitch Horse	0%	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
Go to Shed	0%	20%	40%	20%	<b>80%</b>
Peek out of Cover	60%	<b>100%</b>	80%	<b>100%</b>	<b>100%</b>
Switch Weapon	0%	80%	60%	80%	<b>100%</b>
Combat	0%	0%	0%	0%	<b>20%</b>

# Looking Ahead

## ➤ Modular growth

- ❑ As base model's capabilities improves, Cradle improves accordingly!
- ❑ Some or even all of modules can be combined into a single request.

## ➤ Self-improvement

- ❑ Greatly extends the reach of AI agents. **Every software is potential out-of-box testbed and benchmark.**
- ❑ Cradle collects both high-level semantic actions and low-level keyboard and mouse control for self-improvement.

