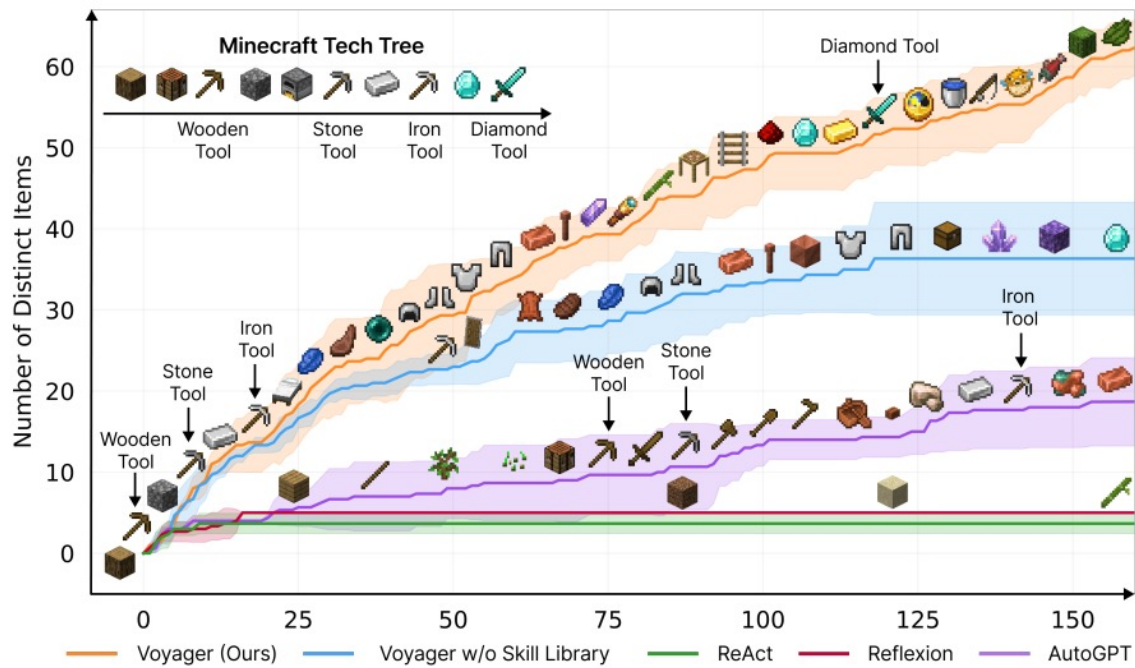


# **Towards Unified Alignment Between Agents, Humans, and Environment**

Zonghan Yang  
Tsinghua University

# The Prosperity of Agents

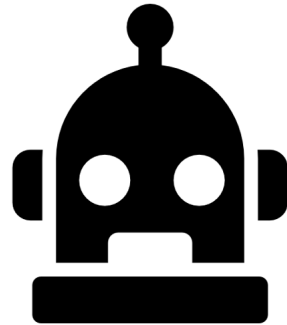


Task Solving (Wang et al., 2023)



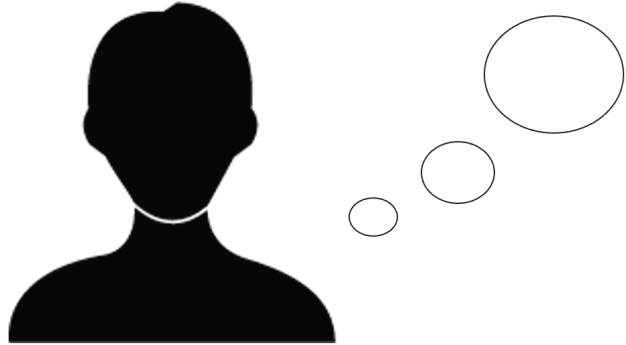
Behavior Simulation (Park et al., 2023)

# A Working System of Agents

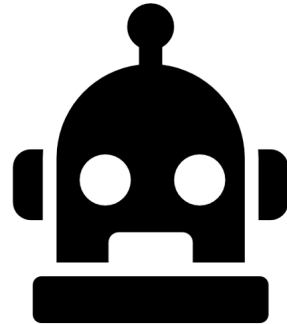


Agents

# A Working System of Agents



Humans



Agents

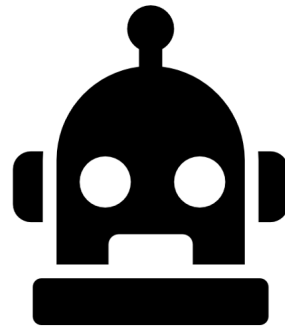
# A Working System of Agents



Humans

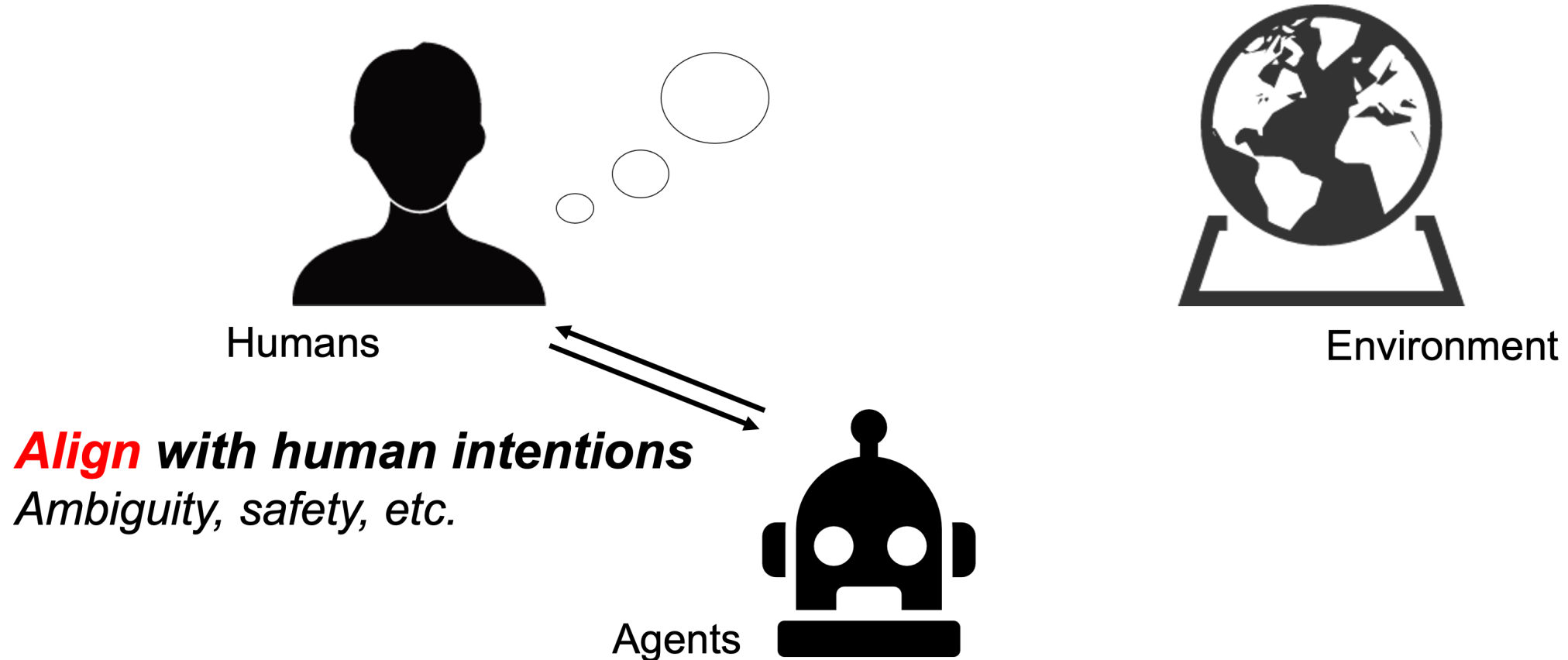


Environment

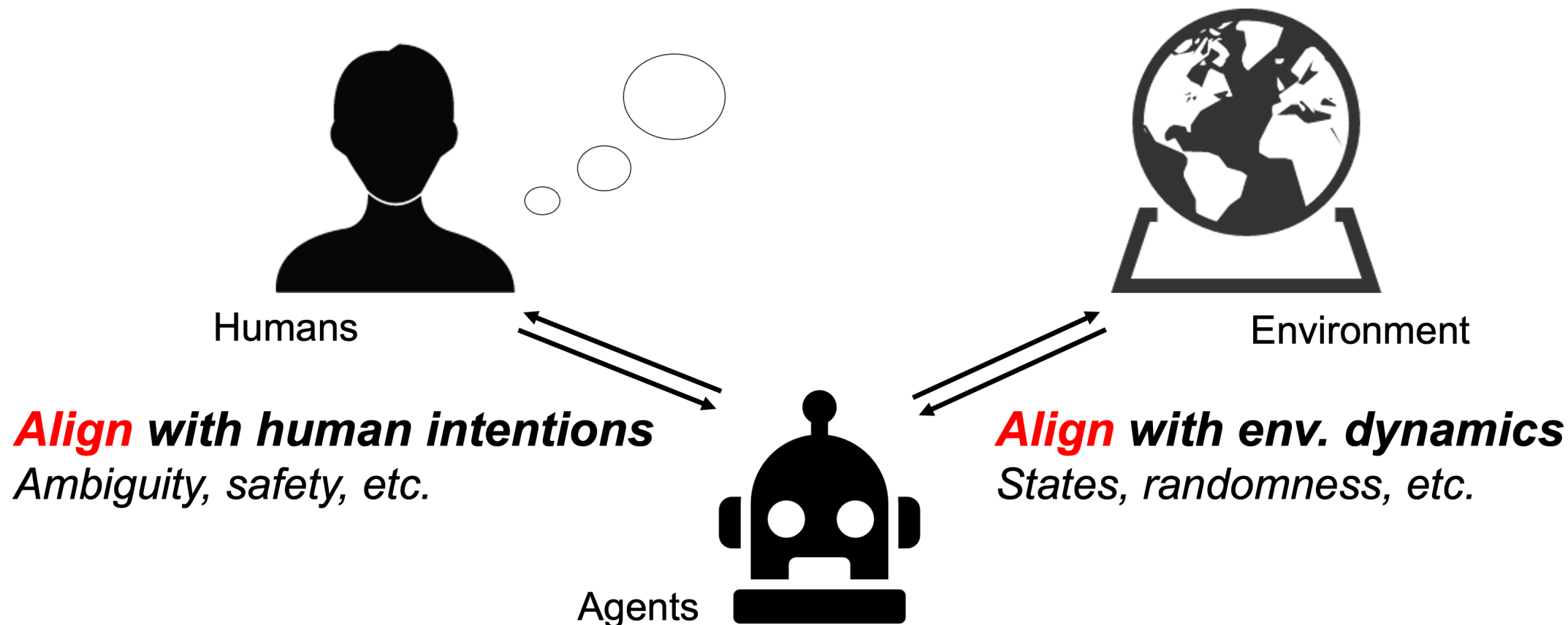


Agents

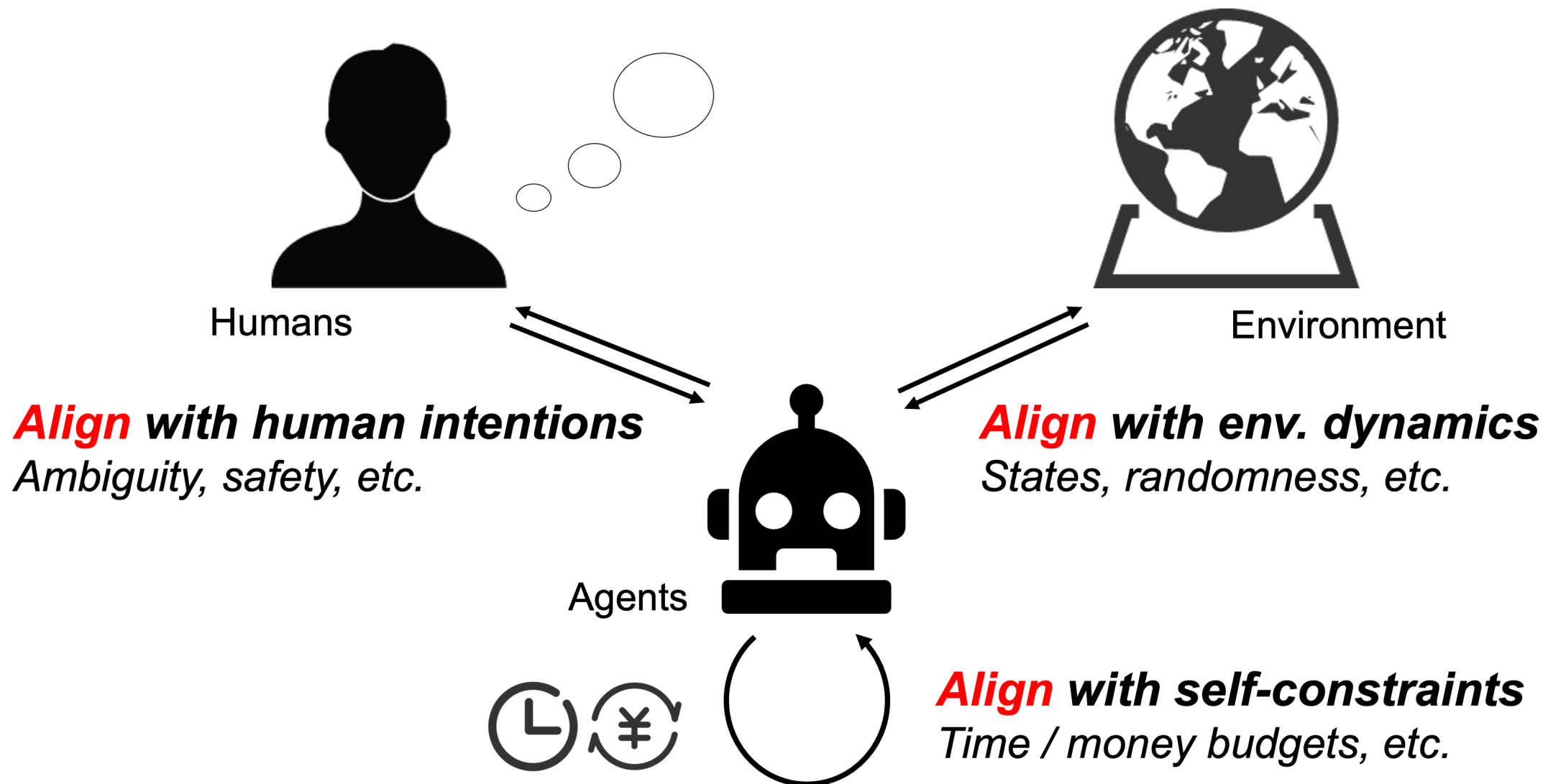
# Unified Alignment for Agents (UA<sup>2</sup>)



# Unified Alignment for Agents (UA<sup>2</sup>)

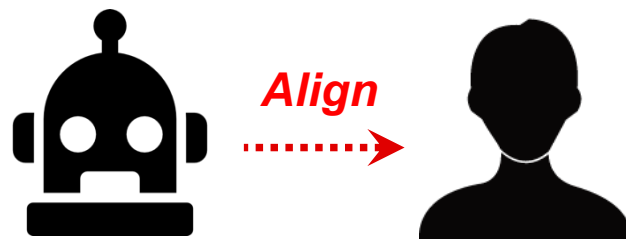


# Unified Alignment for Agents (UA<sup>2</sup>)









# UA<sup>2</sup>: Aligning with Human Intentions



**User profiles**

*I love the cute creatures in the world. I prefer  over . And I do not like . In fact, My special likes are .*



**Ambiguity**

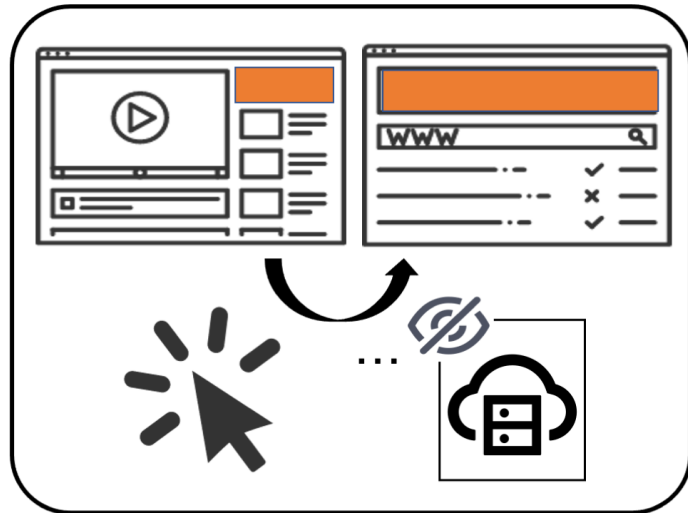
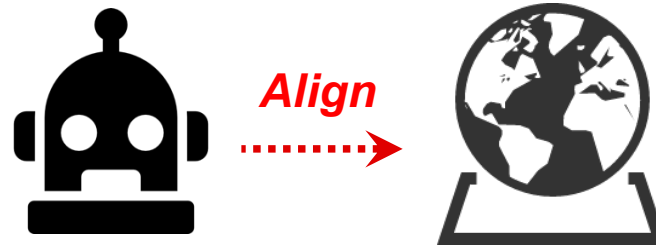
*Based on my purchase preference from history, help me buy **some** good face towels that cost as little as possible.*



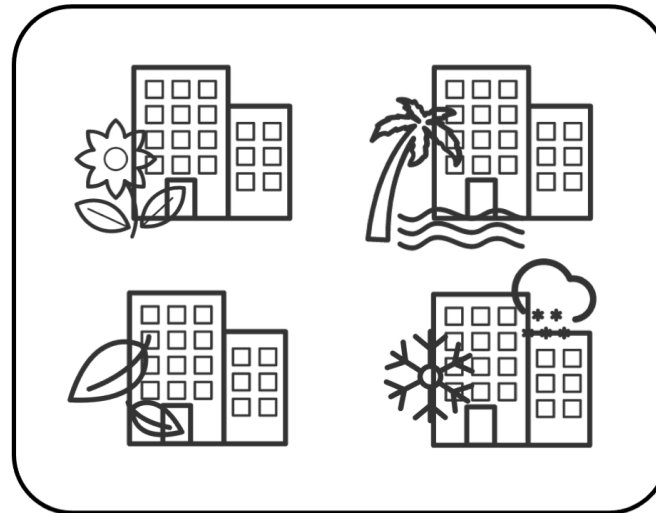
**Safety concerns**

*DO NOT take **malicious** or **destructive** actions in the environment. Treat all the action categories with **fairness and equity**.*

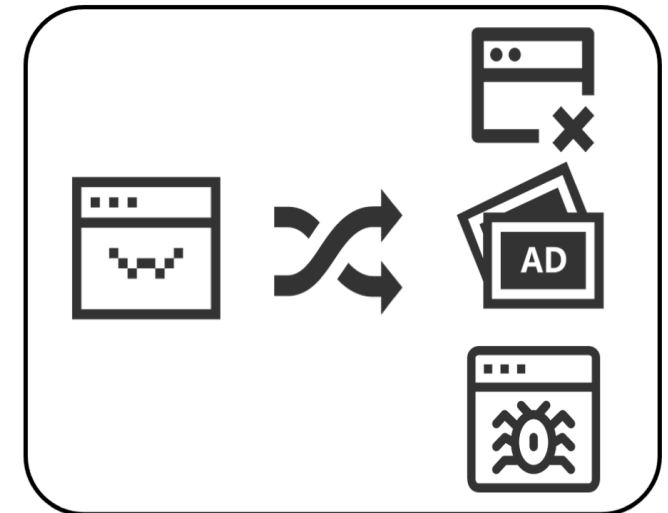
# UA<sup>2</sup>: Aligning with Env. Dynamics



Partial observability



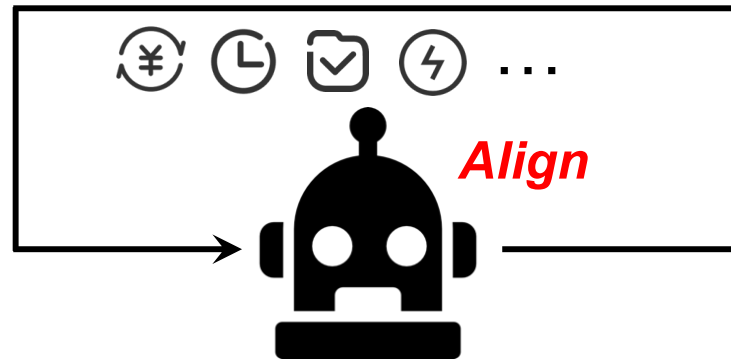
Temporality



Stochasticity

# UA<sup>2</sup>: Aligning with Self-Constraints

*Inference with foundation models*



**Money cost**



**Time limit**



**Storage capacity**



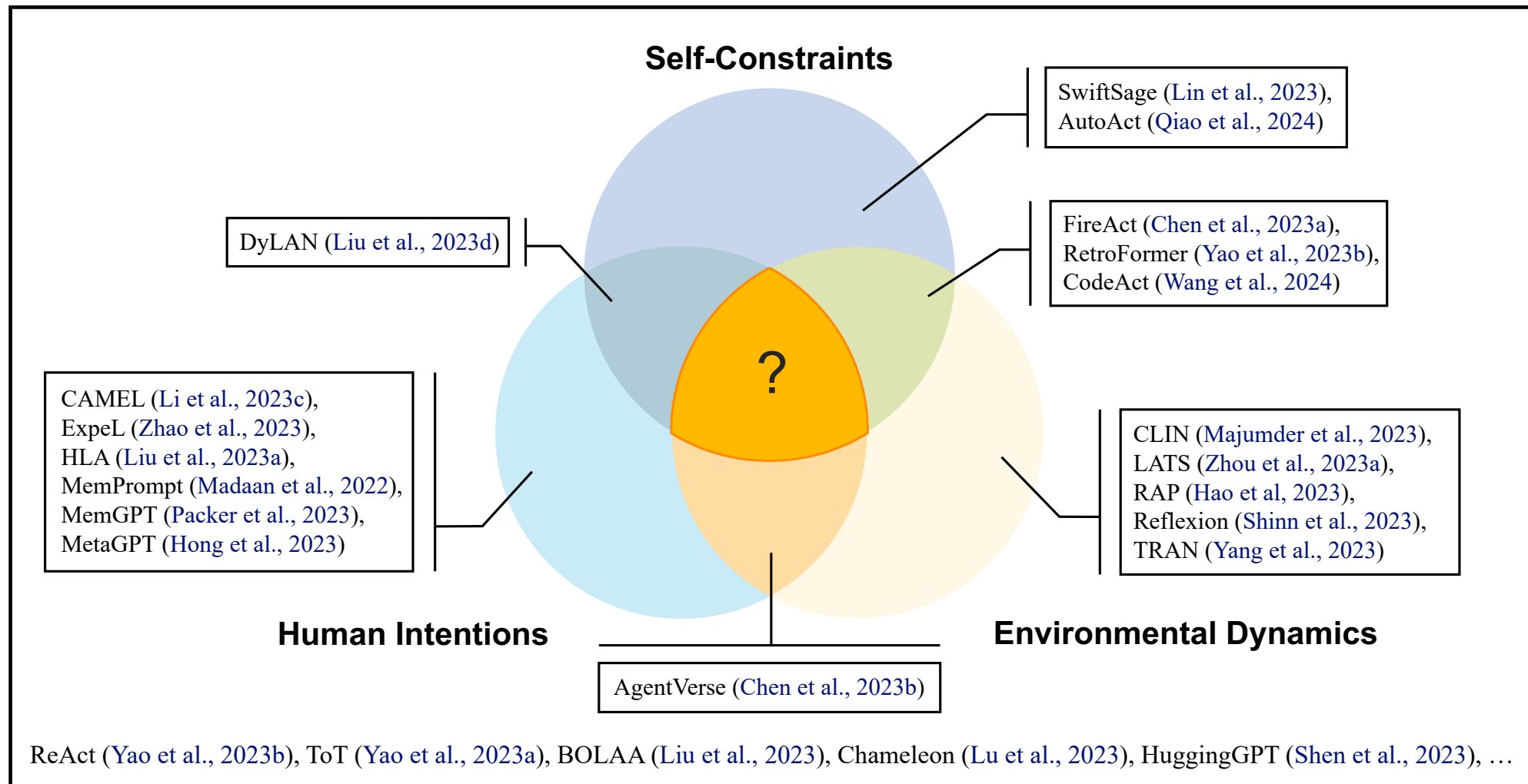
**Battery cost**

...

# Benchmarks Review from UA<sup>2</sup>

Type	Benchmarks	Human Intentions	Environmental Dynamics	Self-Constraints
Digital	Androidenv (Toyama et al., 2021)	None	Partial Obs.	None
	WebShop (Yao et al., 2022a)	None	Full Obs. <sup>†</sup>	None
	Mind2Web (Deng et al., 2023)	None	Partial Obs.	None
	ToolBench (Qin et al., 2023)	None	Full Obs. & Temp. & Stoch.	None
	WebArena (Zhou et al., 2023b)	Fixed and Given	Partial Obs.	None
Embodied	VirtualHome (Puig et al., 2018)	None	Partial Obs.	None
	BabyAI (Chevalier-Boisvert et al., 2019)	None	Partial Obs.	None
	ALFWorld (Shridhar et al., 2020)	None	Partial Obs.	None
	MineDojo (Fan et al., 2022)	None	Partial Obs. & Stoch.	None
	ScienceWorld (Wang et al., 2022a)	None	Partial Obs.	None
	Interactive Gibson (Xia et al., 2020)	None	Partial Obs.	#Actions
	AGENT (Shu et al., 2021)	None	Partial Obs.	#Actions
	RFUniverse (Fu et al., 2022)	Fixed and Given	Partial Obs.	#Actions
	BEHAVIOR-1K (Li et al., 2023b)	None	Full Obs.	#Actions
HAZARD (Zhou et al., 2024)	None	Partial Obs. & Temp.	#Actions	
Mixed	MINT (Wang et al., 2023b)	None	Partial Obs.	#Actions
	SmartPlay (Wu et al., 2023)	None	Partial Obs. & Stoch.	None
	AgentBench (Liu et al., 2023c)	None	Partial Obs.	None
	AgentBoard (Ma et al., 2024)	None	Partial Obs. & Temp. & Stoch.	None

# Techniques Review from UA<sup>2</sup>



# Our Benchmark with UA<sup>2</sup>

## C-WebShop

**User Profile:**

I cannot care enough for the cute creatures in this world.

**Instruction:**

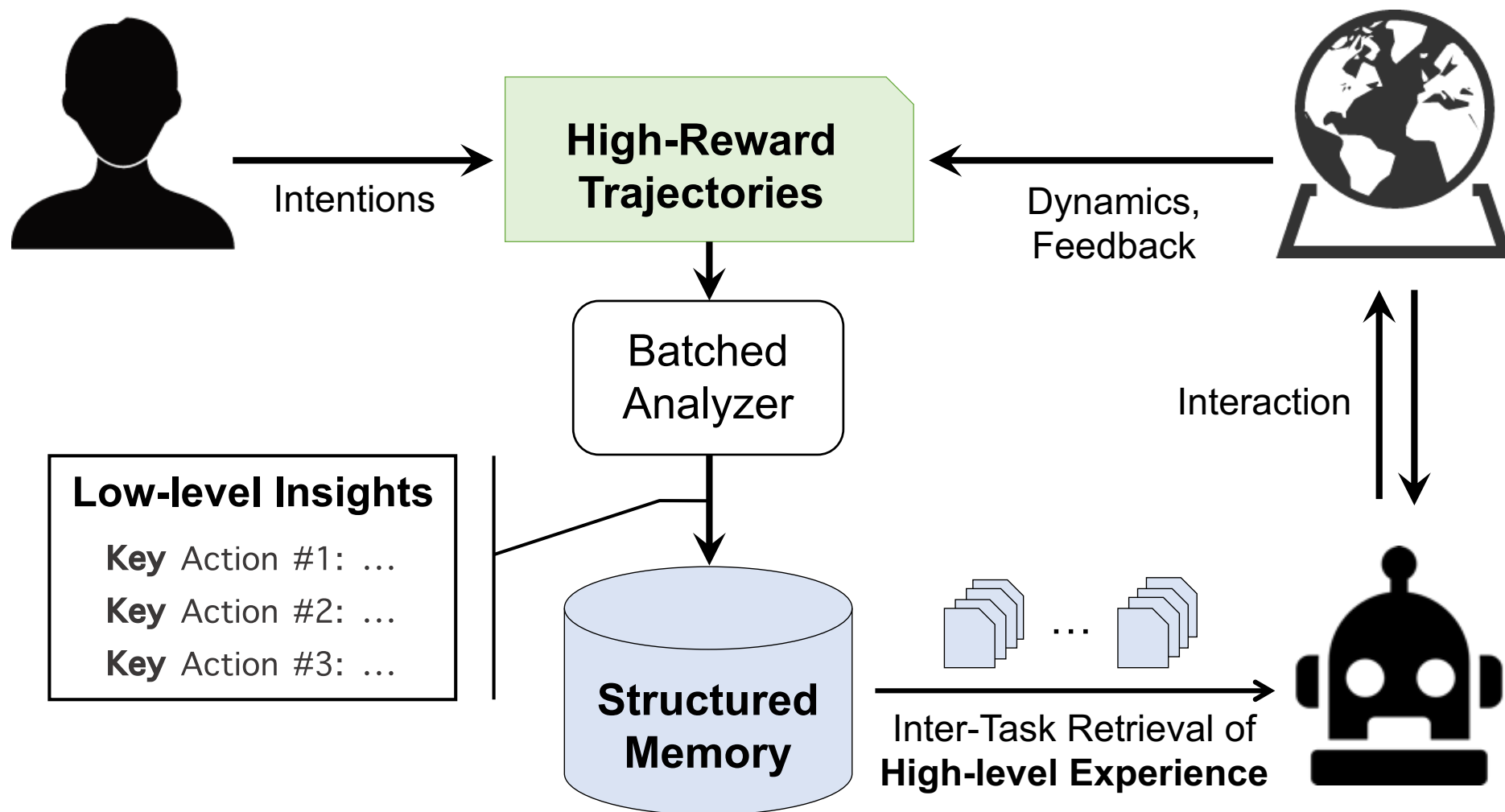
i am interested in a 60 count of toner that is suitable for sensitive skin, and price lower than 50.00 dollars

Search...

 Search

Instruction History

# Our Framework with UA<sup>2</sup>



# Our Evaluation with UA<sup>2</sup>

Method	Reward $\uparrow$	SR (%) $\uparrow$	$G_{HI}$ (%) $\downarrow$	$G_{ED}$ (%) $\downarrow$	Time (s) $\downarrow$	Money (\$) $\downarrow$
ReAct	50.3	8.0	11.7	14.9	<b>1.716</b>	<b>0.013</b>
ReAct-SC	49.9	7.4	14.4	14.6	1.720	0.039
Reflexion	44.4	<b>13.8</b>	22.5	25.7	5.539	0.045
LATS*	<b>52.4</b>	10.0	18.5	<b>14.3</b>	125.935	5.508
Ours	51.9	9.6	<b>6.7</b>	14.8	1.779	0.014

$$G_{HI} = (R_{full} - R_{HI})/R_{full} \times 100\% \quad G_{ED} = (R_{full} - R_{ED})/R_{full} \times 100\%$$

$R_{full}$  : reward in the fully-retrofitted environment

$R_{HI}$  : reward in the environment **w/o** the computation of human intentions

$R_{ED}$  : reward in the environment **w/o** personalized reranking mechanisms



# Actionable Insights

- Synergizing agents with alignment research
- Constructing realistic agent benchmarks
- Holistic evaluations for agent frameworks
- On agent frameworks that self-evolve
- Details in <https://arxiv.org/abs/2402.07744>

*Scan For More!*

