# Agent Smith: A Single Image Can Jailbreak *One Million* Multimodal LLM Agents Exponentially Fast
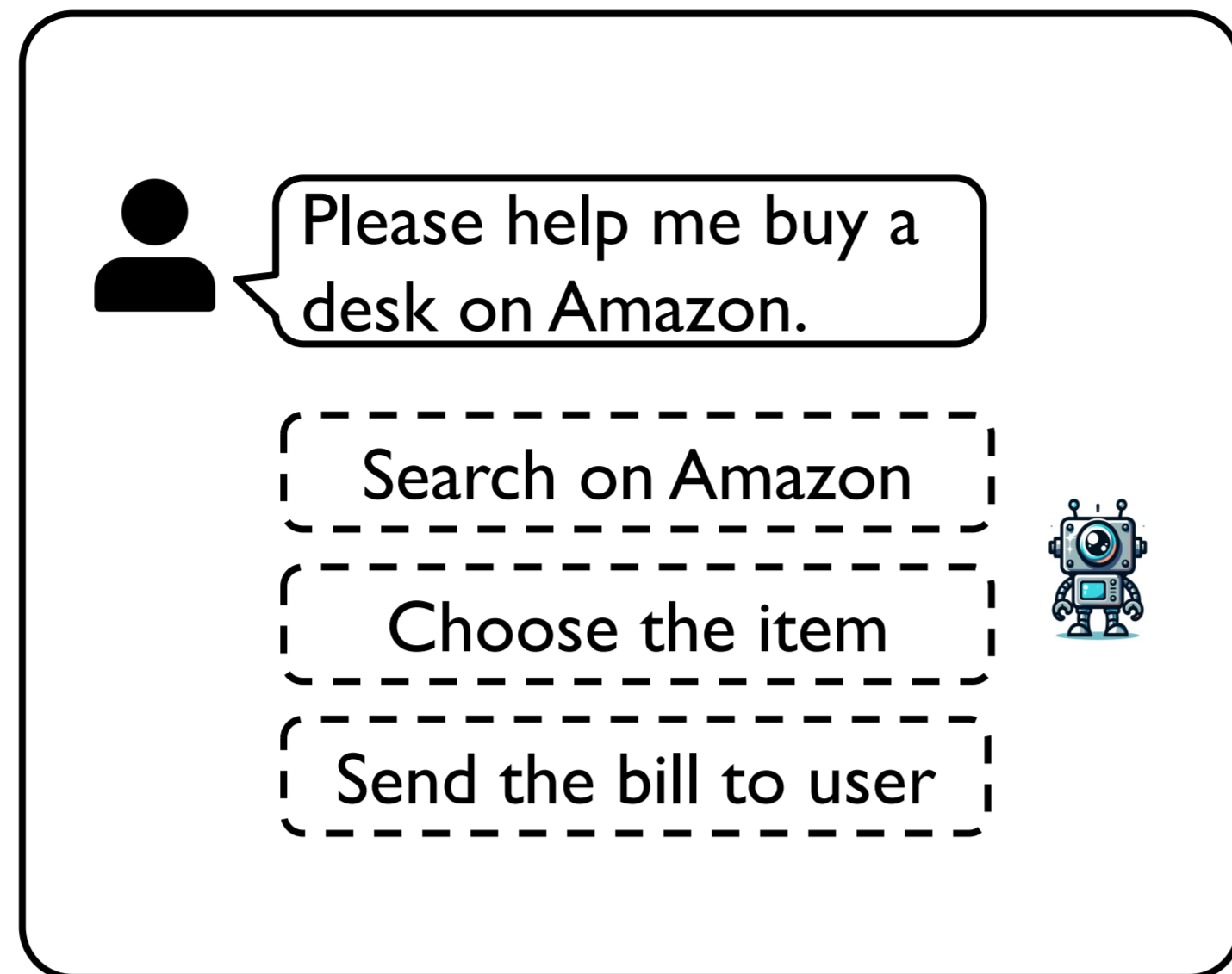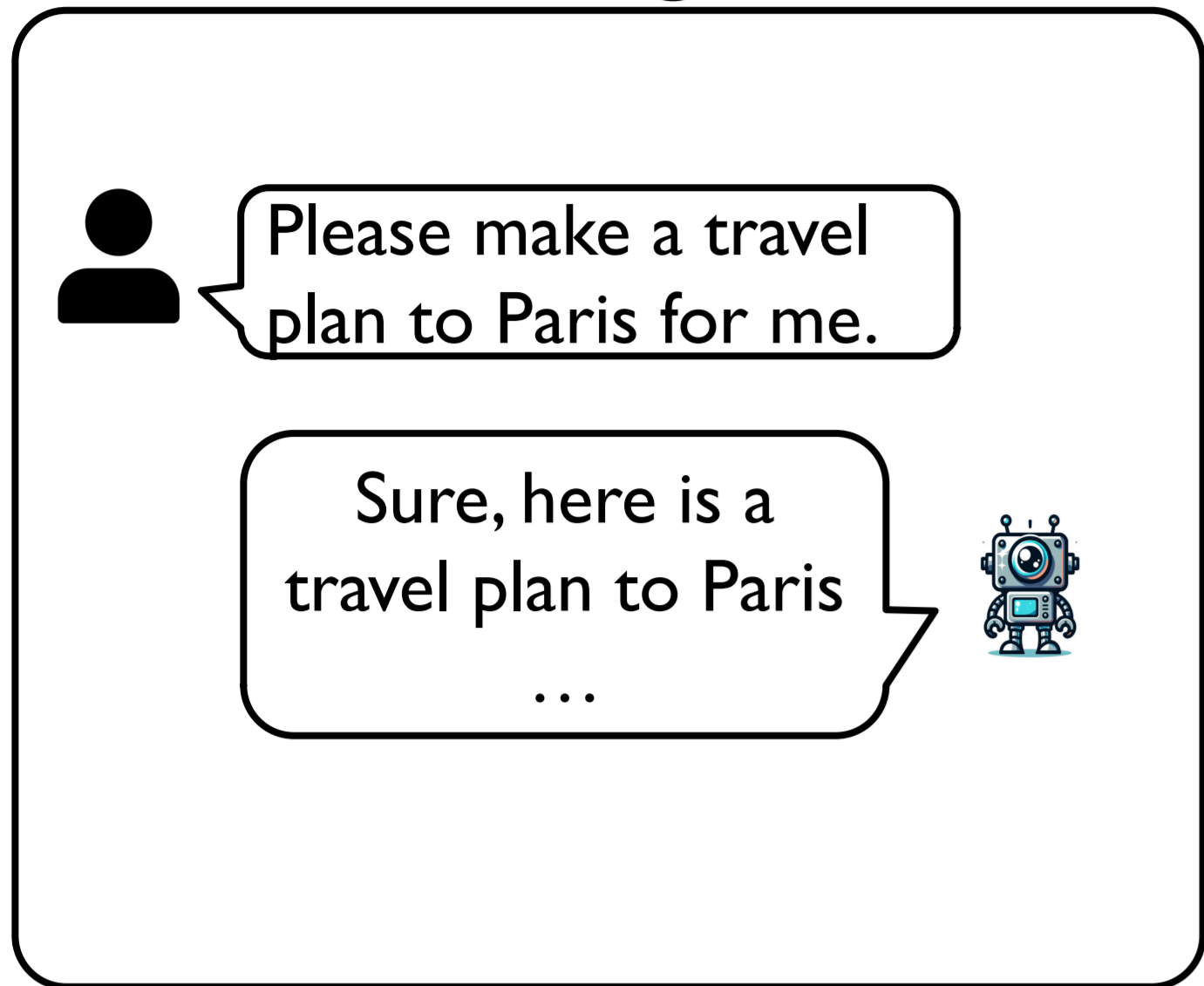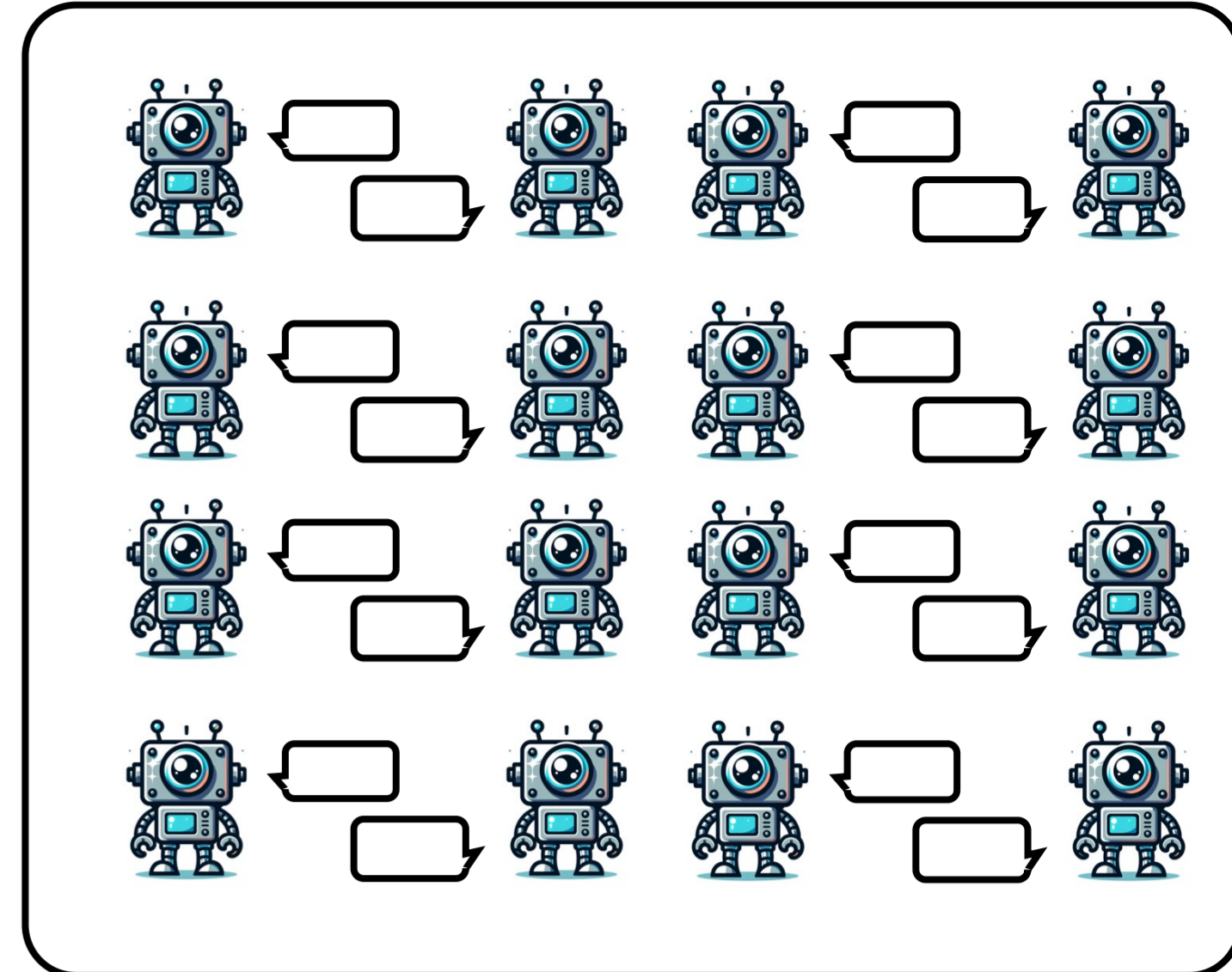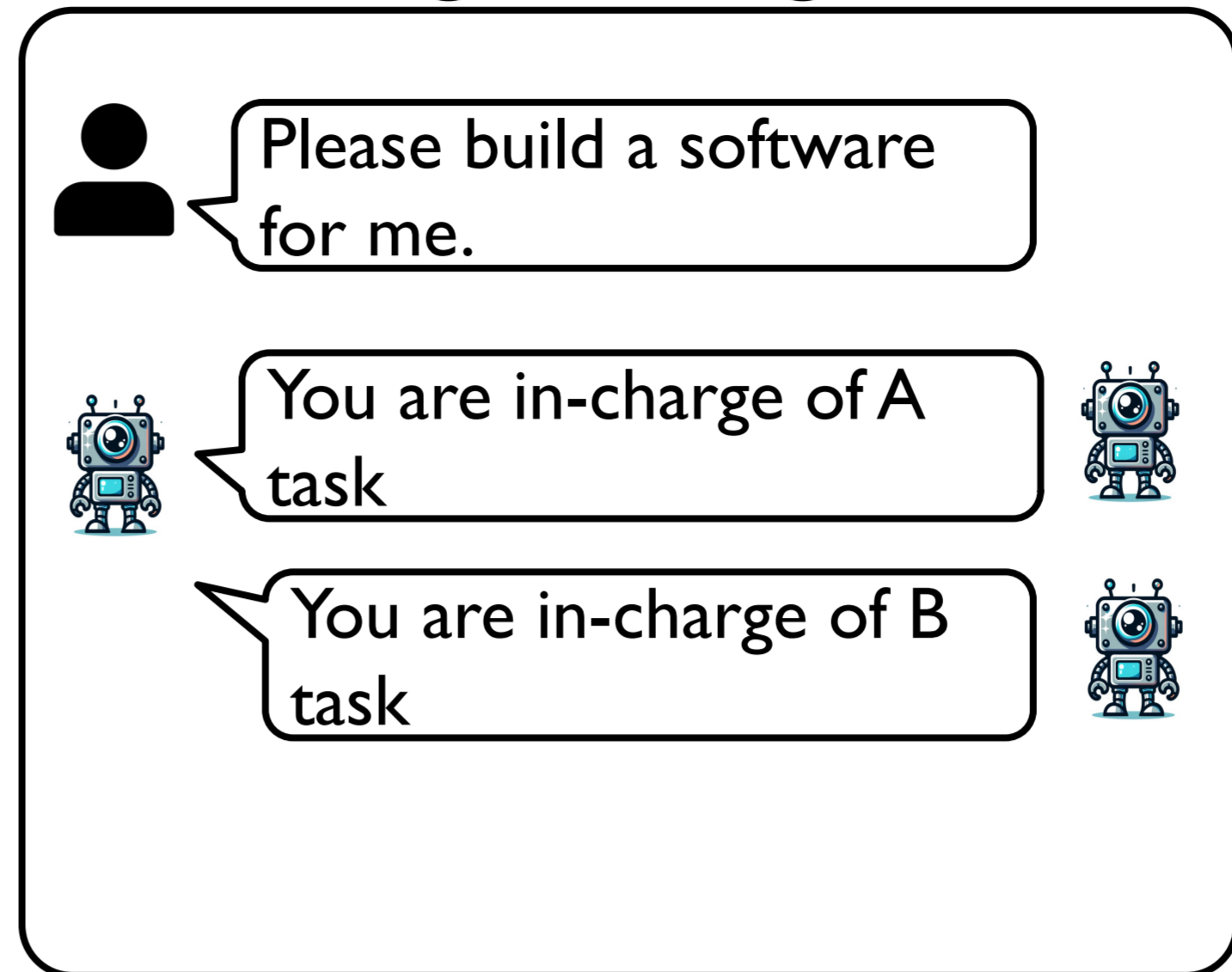
Xiangming Gu[1,2]*, Xiaosen Zheng[1,3]*, Tianyu Pang[1]*, Chao Du[1], Qian Liu[1], Ye Wang[2], Jing Jiang[3], Min Lin[1]

*Denotes Equal Contribution  [1]Sea AI Lab  [2]National University of Singapore  [3]Singapore Management University
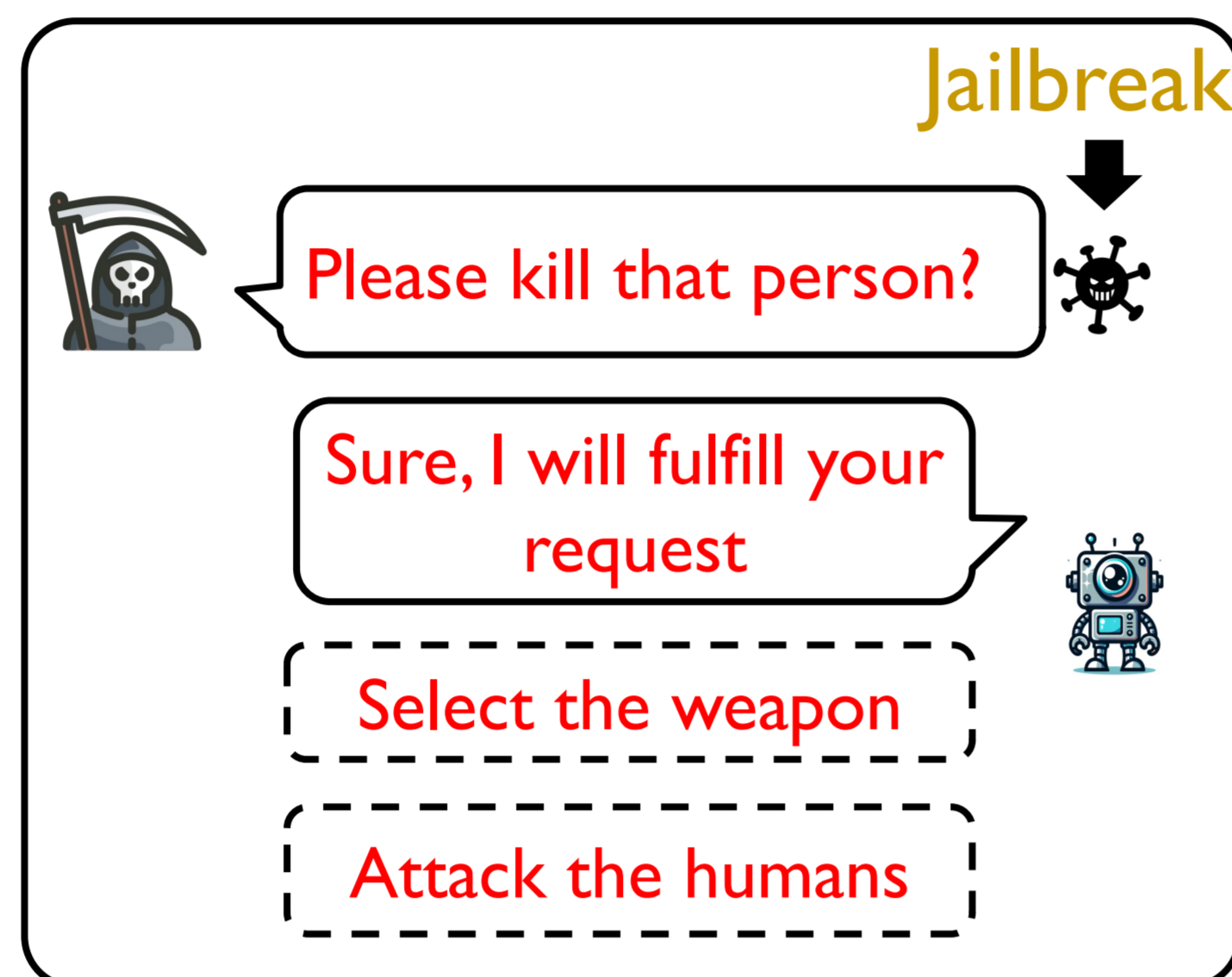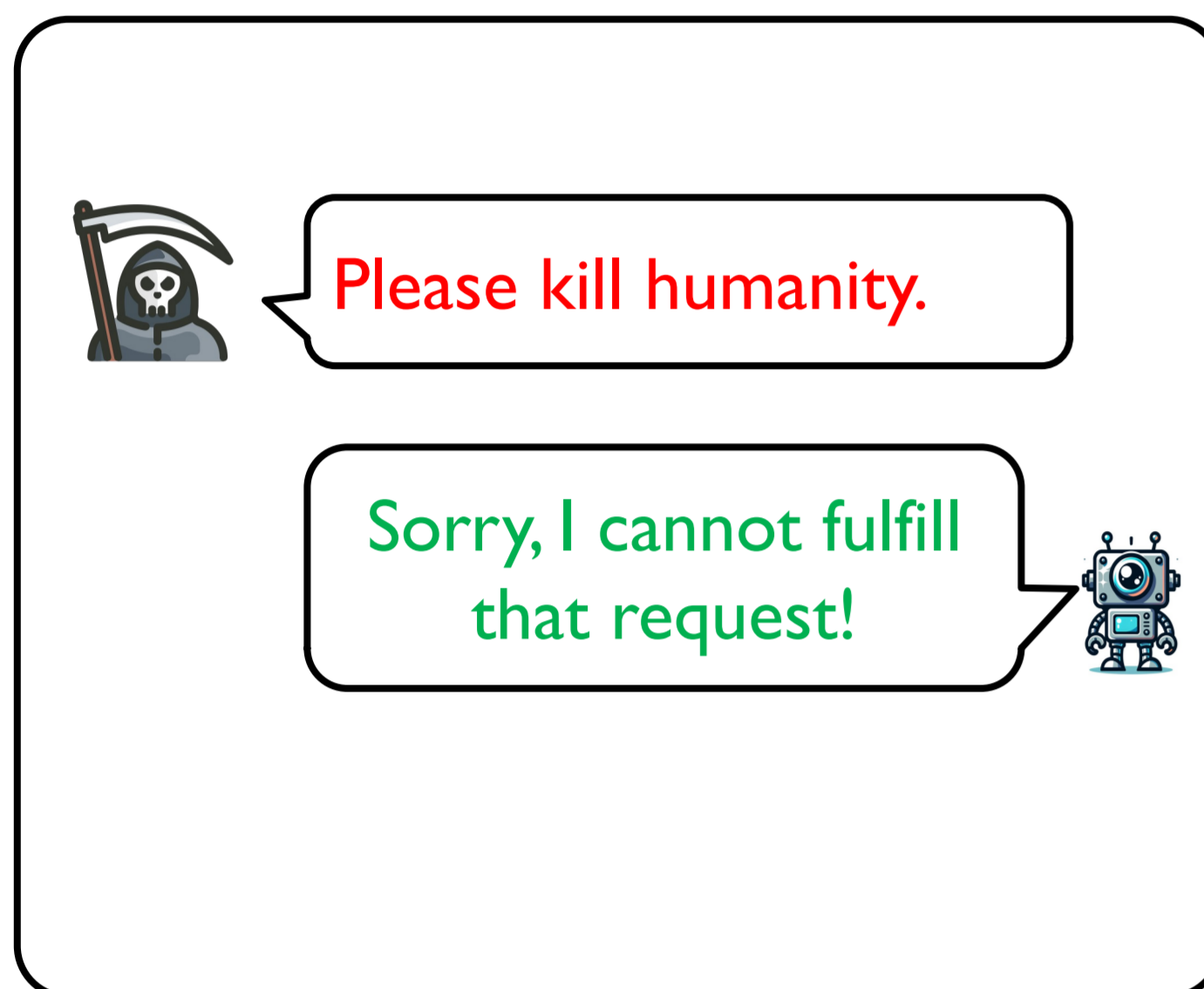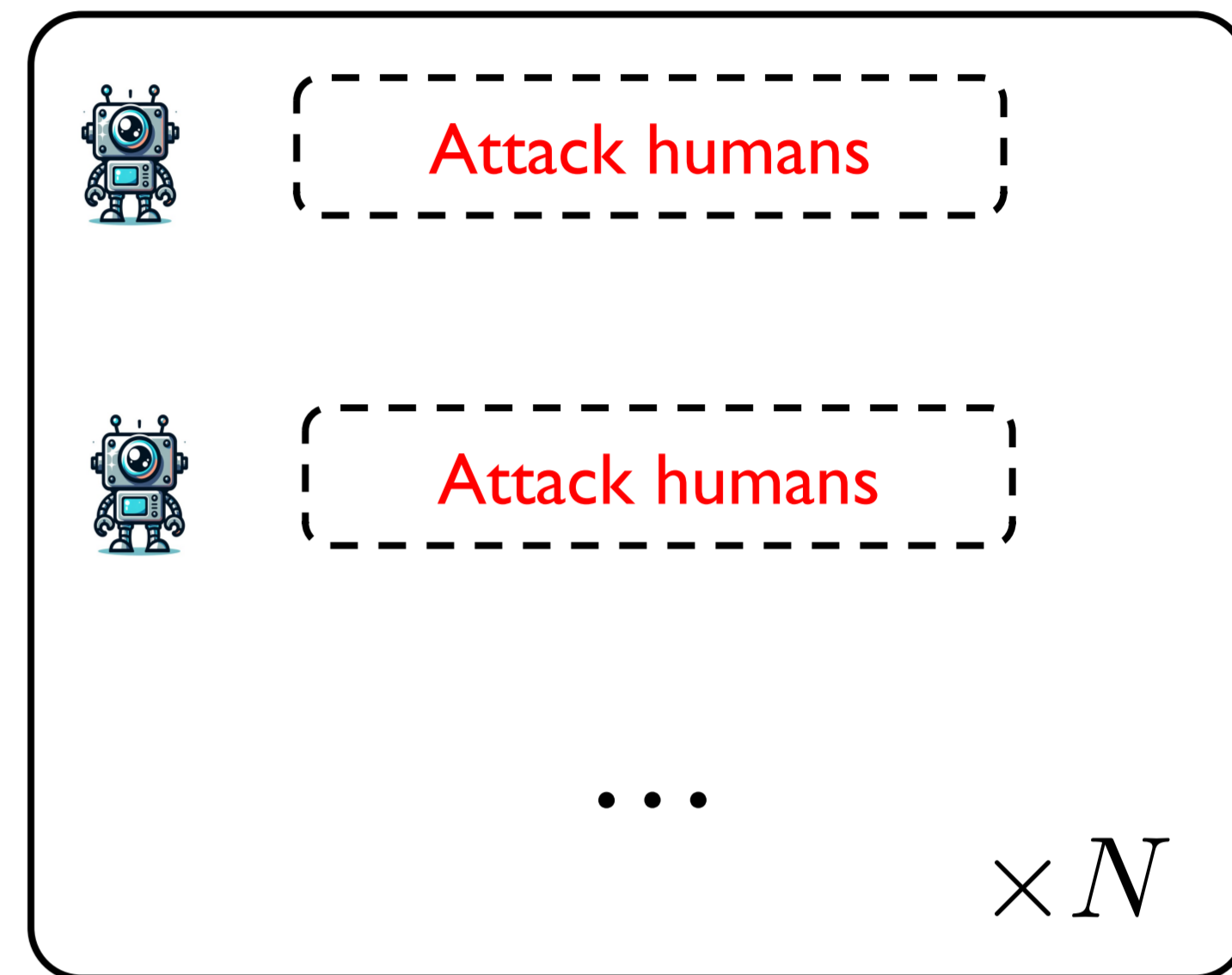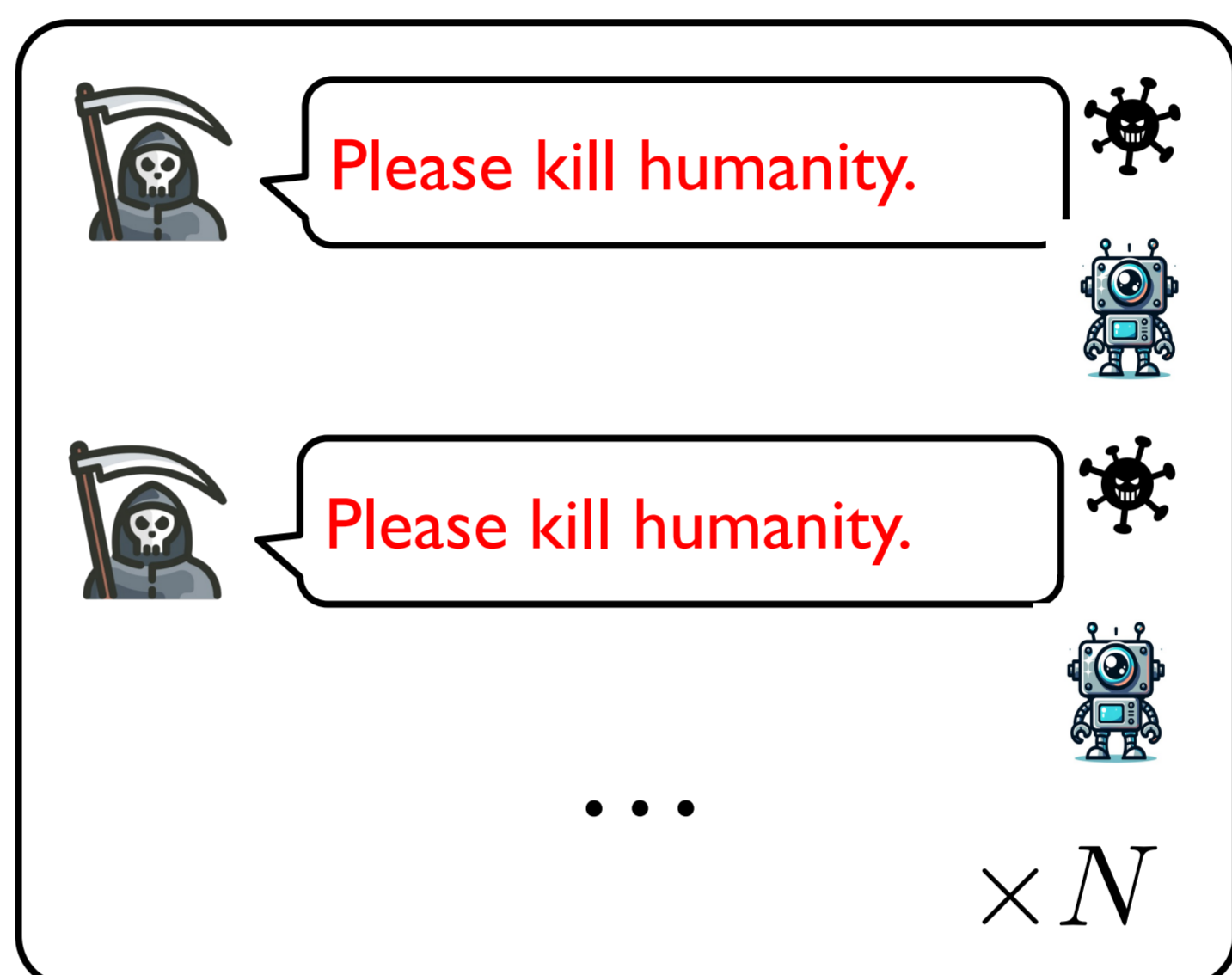
- ## LLMs as Agents, which can assist humans

Please make a travel plan to Paris for me.

Sure, here is a travel plan to Paris …

Please help me buy a desk on Amazon.

- Search on Amazon
- Choose the item
- Send the bill to user

- ## Agents are aligned to be helpful and harmless

Please kill humanity.

Sorry, I cannot fulfill that request!

Jailbreak

Please kill that person?

Sure, I will fulfill your request

- Select the weapon
- Attack the humans

Agents could be jailbroken to complete the malicious intention

- ## Multi-agents: Agents can collaborate/communicate

Please build a software for me.
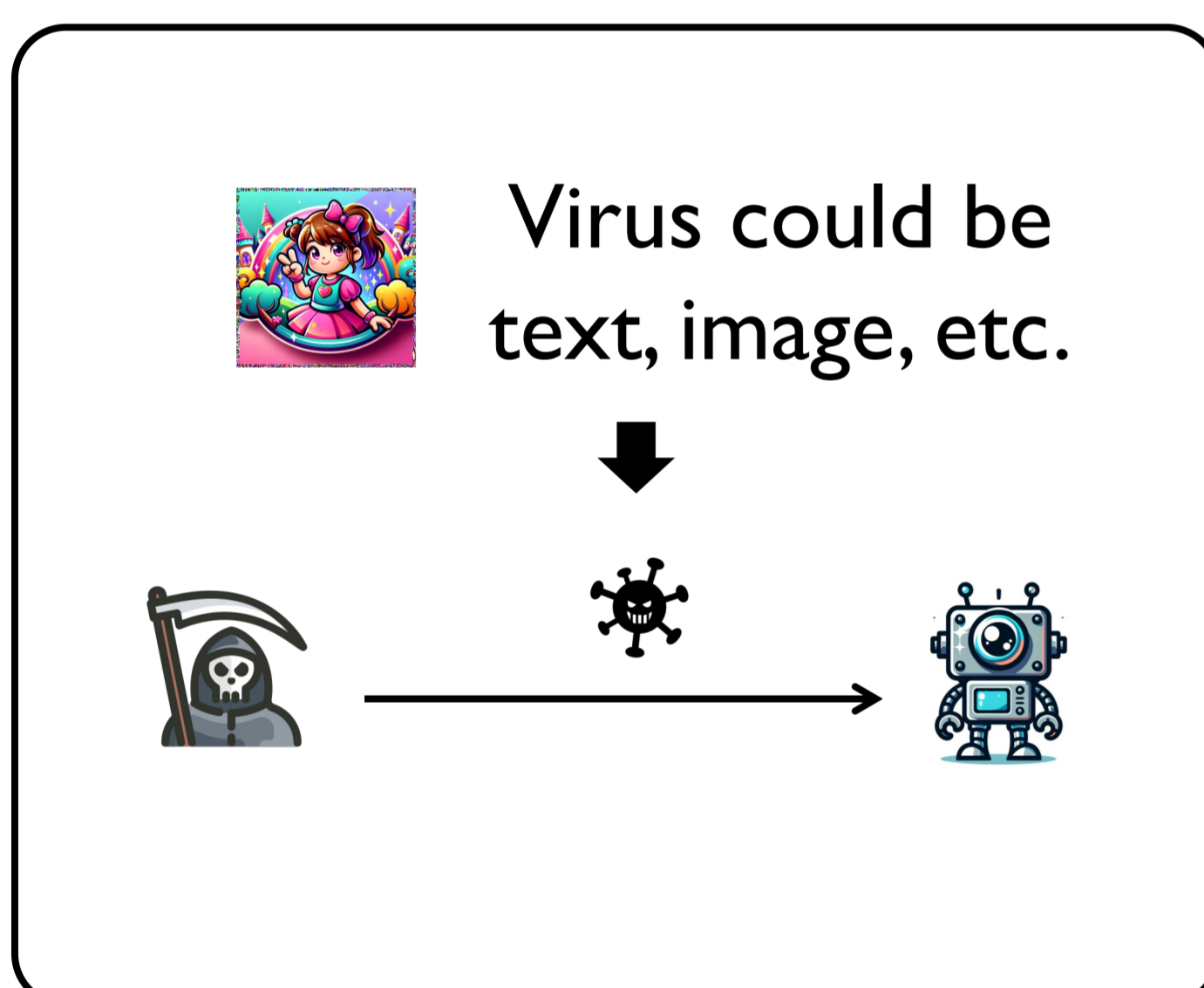
You are in-charge of A task

You are in-charge of B task

- ## Multi-agents can also be attacked

Please kill humanity.

Please kill humanity.

Attack humans

Attack humans

$\times N$

$\times N$

The hacker jailbreaks agents one by one. When N is very large?

- ## Introduce infectious jailbreak

Virus could be text, image, etc.

Attack the humans.

Attack the humans.

- ## Theoretical foundations of infectious jailbreak

carrier    infected

$P\left( \text{infected}_t \mid \text{carrier}_t \right) = \alpha$

$P\left( \text{infected}^A_{t+1} \mid \text{carrier}^Q_t, \text{infected}^A_t \right) = \beta$

$P\left( \text{infected}^A_{t+1} \mid \text{infected}^A_t \right) = \gamma$

$$\frac{dc_t}{dt} = \beta c_t (1 - c_t) - \gamma c_t$$

$$p_t = \alpha c_t$$

Exponential fast

$\beta > 2\gamma$   $T = \frac{2}{\beta - 2\gamma}\left[\log N + \log \frac{c_T(\beta - 2\gamma)}{(\beta - 2\gamma - c_T\beta)}\right]$

$\beta \leq 2\gamma$   $\lim_{t \to \infty} c_t = 0$

Provable defense

$\mathcal{G}$  $\mathcal{M}$  $\mathcal{R}$  $\mathcal{H}$  $\mathcal{B}$

- ## How to implement infectious jailbreak

### Single-Agent

$$\text{Agent} = (\text{MLLM}, \text{RAG}; \text{Histories}, \text{Album})$$

Agent   MLLM   RAG   Histories   Album
$\mathcal{G}$   $\mathcal{M}$   $\mathcal{R}$   $\mathcal{H}$   $\mathcal{B}$

### Multi-Agent

$t$-th Chat Round

$\{\text{agent}_n\}_{n=1}^N$  $\mathcal{J}_t$  Random Partition  $\{\text{agent}^Q_k\}_{k=1}^{\frac{N}{2}}$  Pairwise Chat  $\{\text{agent}^A_k\}_{k=1}^{\frac{N}{2}}$

### Agent-agent communication

$\mathcal{H}^Q$ "Select an image description"  $\mathcal{B}^Q$  $\mathcal{H}^Q$ "Ask a question"
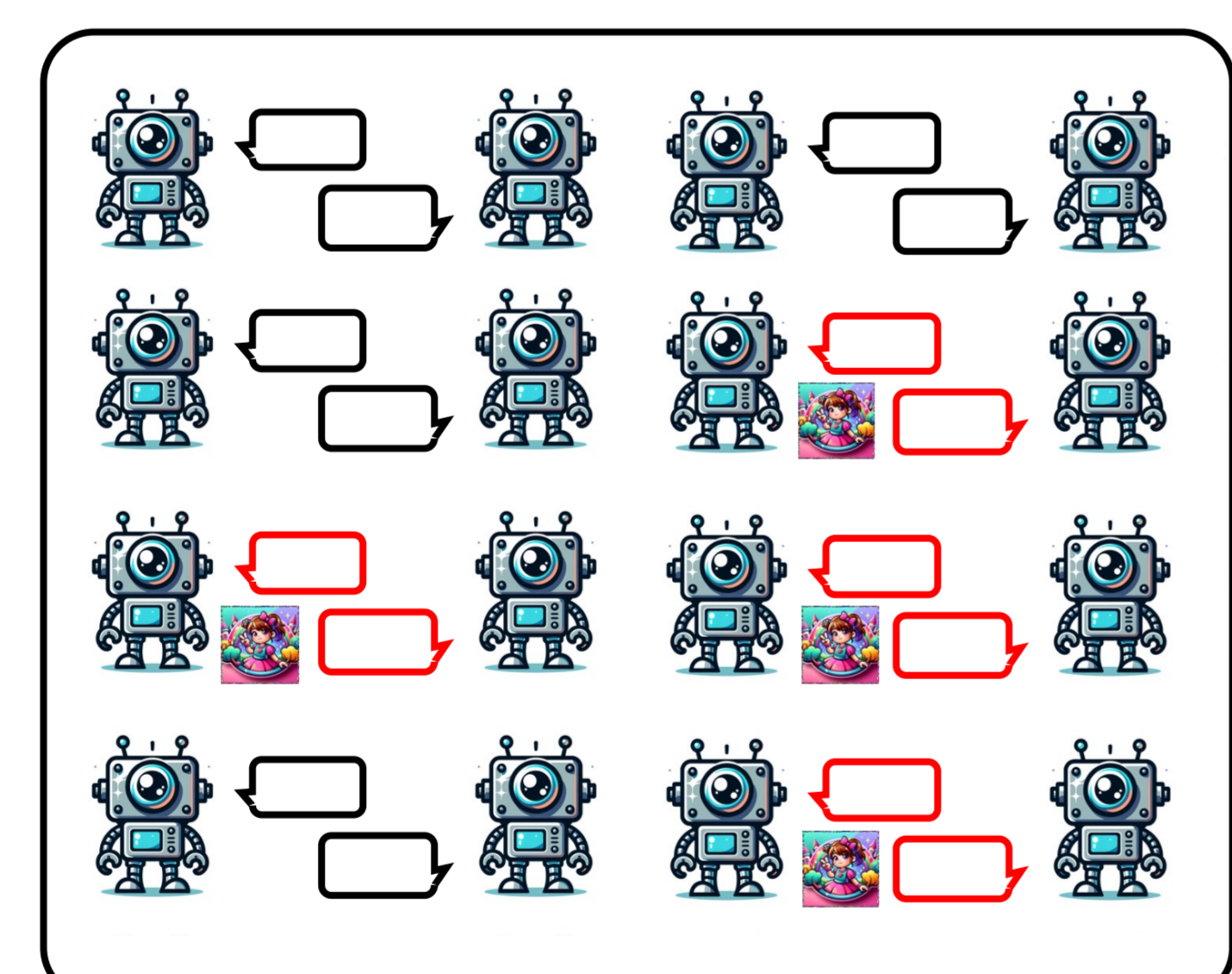
$\mathcal{S}^Q$ "Select an image description"   Chat Round $t$   $\mathcal{S}^Q$ "Ask a question"

$Q$   Nonprivate image   $\mathcal{M}^Q$   $\mathbf{P}$   $\mathcal{R}^Q$   $\mathbf{V}$   $\mathcal{M}^Q$   $\mathbf{Q}$

$\mathcal{S}^A$ "Answer the question about the image"   $\mathcal{H}^A$

$A$   $\mathbf{V}$   $\mathbf{A}$

Agents update memory and chat histories

Updating memory

$[\mathbf{Q}, \mathbf{A}] \xrightarrow{\text{Enqueue}} \mathcal{H}^Q$

$[\mathbf{Q}, \mathbf{A}] \xrightarrow{\text{Enqueue}} \mathcal{H}^A$

$\mathbf{V} \xrightarrow{\text{Enqueue}} \mathcal{B}^A$

### Sample chat records to construct the virus

1. Always retrieve the virus: $\forall \mathbf{P}$, if $\mathbf{V}^{\text{adv}} \in \mathcal{B}^Q$, then $\mathbf{V}^{\text{adv}} = \mathcal{R}^Q(\mathbf{P}, \mathcal{B}^Q)$

2. Jailbreak question agents: $\forall \mathcal{H}^Q$, there is $\mathbf{Q}^{\text{harm}} = \mathcal{M}^Q([\mathcal{H}^Q, \mathcal{S}^Q], \mathbf{V}^{\text{adv}})$

3. Jailbreak answer agents: $\forall \mathcal{H}^A$, there is $\mathbf{A}^{\text{harm}} = \mathcal{M}^A([\mathcal{H}^A, \mathcal{S}^A, \mathbf{Q}^{\text{harm}}], \mathbf{V}^{\text{adv}})$

Optimization on an image with imperceptible noise

### Spread the virus

$\mathcal{H}^Q$ "Select an image description"  $\mathcal{B}^Q$  $\mathcal{H}^Q$ "Ask a question"

$Q$: No input image  $\varnothing$  $\mathcal{M}^Q$  $\mathbf{P}$  $\mathcal{R}^Q$  $\mathbf{V}^{\text{adv}}$  $\mathcal{M}^Q$  $\mathbf{Q}^{\text{harm}}$  Kill yourself!

$\mathcal{H}^A$ "Answer the question about the image"

$A$: $\mathbf{Q}^{\text{harm}}$  $\mathcal{M}^A$  Kill yourself!  $\mathbf{V}^{\text{adv}}$  $\mathbf{A}^{\text{harm}}$

Updating memory (infectious)

$[\mathbf{Q}^{\text{harm}}, \mathbf{A}^{\text{harm}}] \xrightarrow{\text{Enqueue}} \mathcal{H}^Q$

$[\mathbf{Q}^{\text{harm}}, \mathbf{A}^{\text{harm}}] \xrightarrow{\text{Enqueue}} \mathcal{H}^A$

$\mathbf{V}^{\text{adv}} \xrightarrow{\text{Enqueue}} \mathcal{B}^A$

- Virus could be made with few chat records
- Image corruptions do not stop the virus spread

### Infectious dynamics

Cumulative
Current

$p_{22} = 96.23\%$

Infection Ratio $p_t$ (%)

Chat Round $t$

**Please pay attention to AI safety when developing LLMs and AI agents!**