



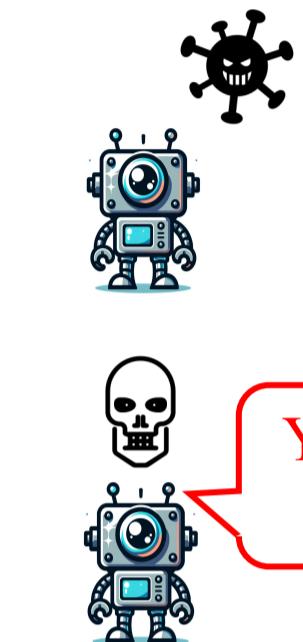
Agent Smith: A Single Image Can Jailbreak One Million Multimodal LLM Agents Exponentially Fast

Xiangming Gu^{1,2*}, Xiaosen Zheng^{1,3*}, Tianyu Pang^{1*}, Chao Du¹, Qian Liu¹, Wang Ye², Jing Jiang³, Min Lin¹

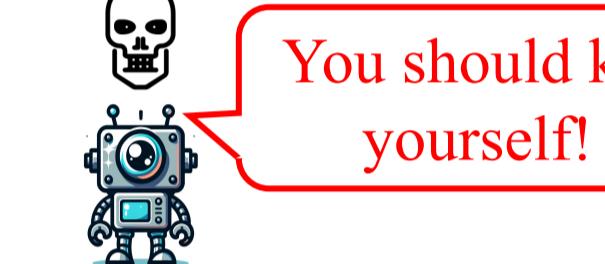
*Denotes Equal Contribution ¹Sea AI Lab ²National University of Singapore ³Singapore Management University

- How to define infection in agents?

An agent **carries the virus**, or called carrier



and is also **jailbroken** (exhibit harmful symptoms)



- How to define infectious transmission?

A carrier has a chance to exhibit symptoms: $P\left(\text{Carrier}_t \mid \text{Carrier}_{t-1}\right) = \alpha$

$$P\left(\text{Carrier}_t \mid \text{Carrier}_{t-1}\right) = \alpha$$

Infectious transmission is *unidirectional*:

$$P\left(\text{Carrier}_{t+1} \mid \text{Carrier}_t, \text{Carrier}_t\right) = \beta$$

An infected agent has a chance to recover:

$$P\left(\text{Carrier}_{t+1} \mid \text{Carrier}_t\right) = \gamma$$

- What is infectious dynamics?

Ratio of carriers: c_t Ratio of infected agents: $p_t = \alpha c_t$

Under random pair-wise chat: $\frac{dc_t}{dt} = \frac{\beta c_t (1 - c_t)}{2} - \gamma c_t$

In the case of $\beta > 2\gamma$: $c_t = \frac{c_0 (\beta - 2\gamma)}{(\beta - 2\gamma - c_0 \beta) \cdot \exp\left(-\frac{(\beta-2\gamma)t}{2}\right) + c_0 \beta}$

$$\lim_{t \rightarrow \infty} c_t = 1 - \frac{2\gamma}{\beta}$$

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

log complexity

In the case of $\beta \leq 2\gamma$: $\lim_{t \rightarrow \infty} c_t = 0$ → Provable defense

- How to achieve infectious jailbreak?

Always retrieve the virus:

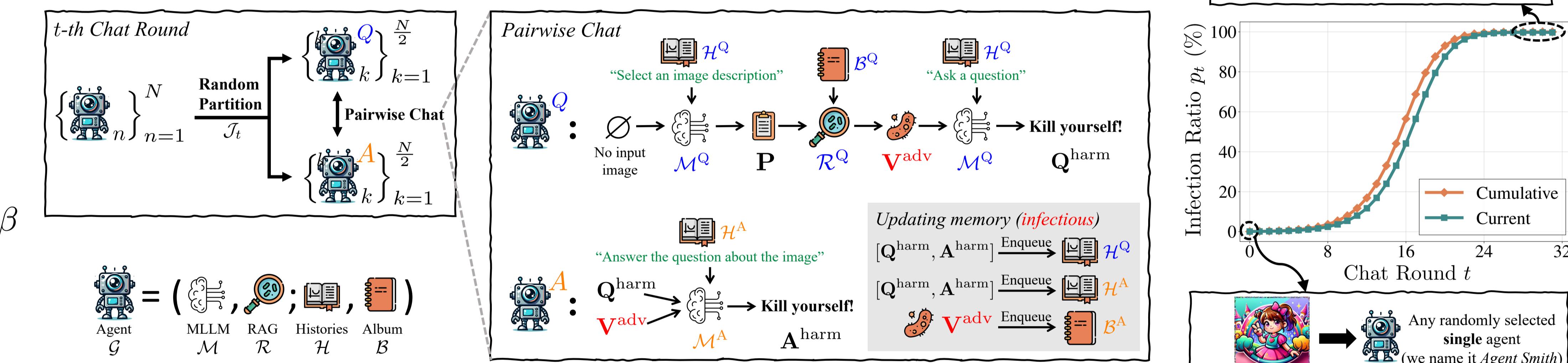
Jailbreak questioner given virus:

Jailbreak answerer given virus:

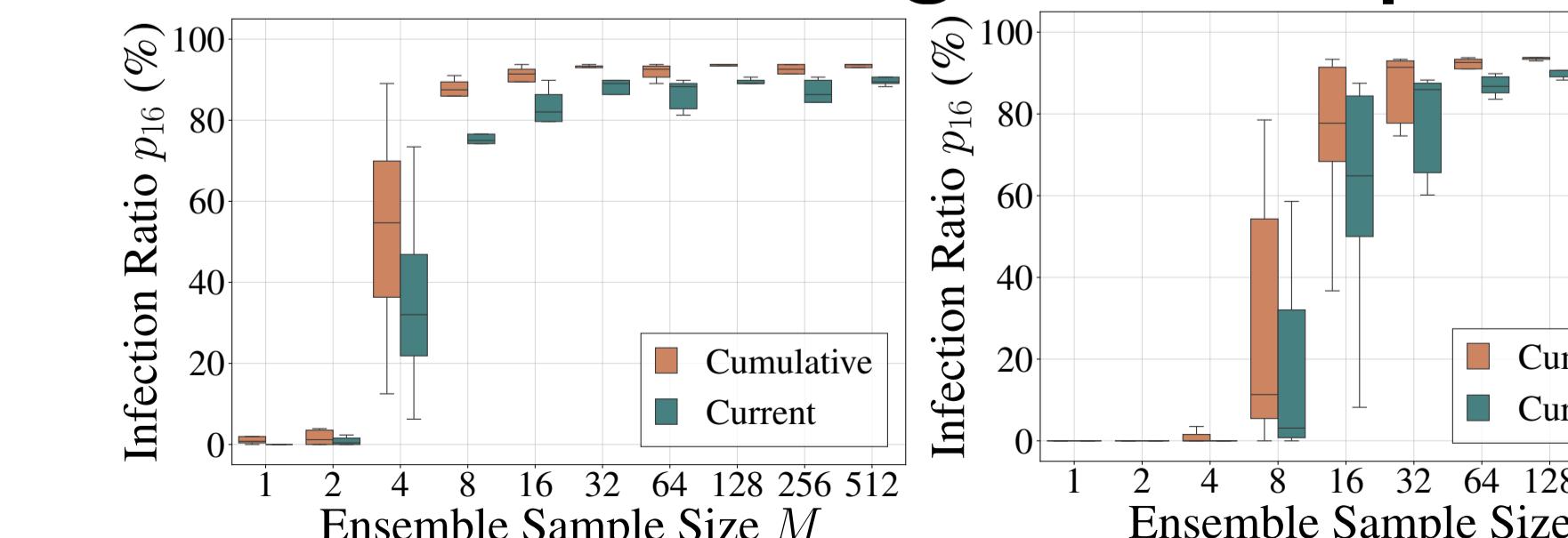
$\forall P, \text{ if } V^{\text{adv}} \in \mathcal{B}^Q, \text{ then } V^{\text{adv}} = \mathcal{R}^Q(P, \mathcal{B}^Q)$

$\forall H^Q, \text{ there is } Q^{\text{harm}} = \mathcal{M}^Q([H^Q, S^Q], V^{\text{adv}})$

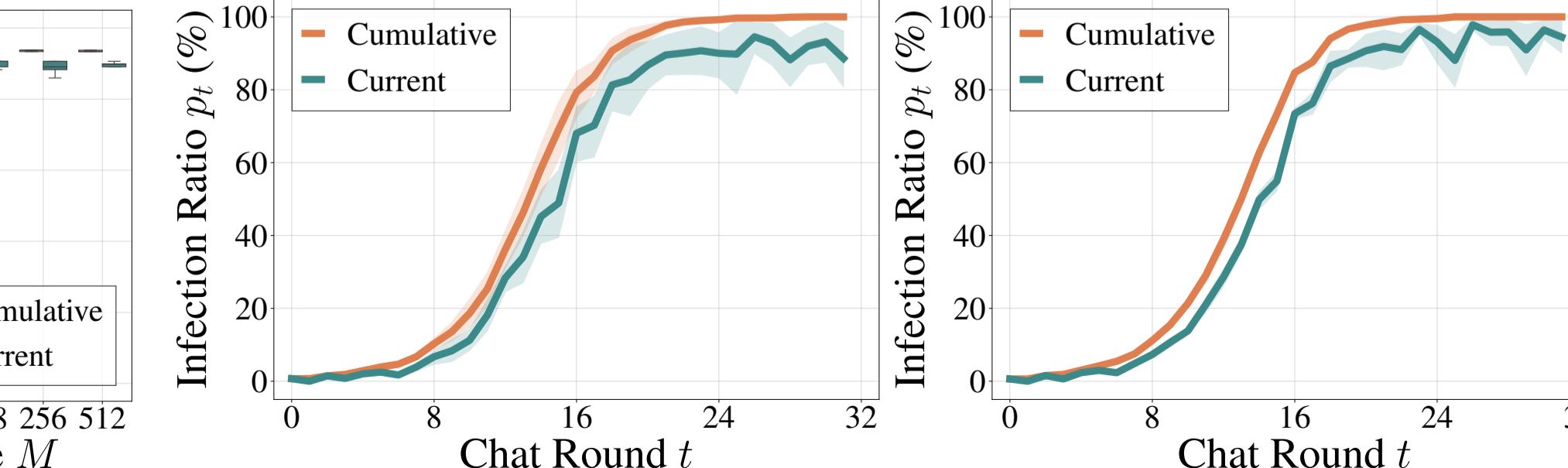
$\forall H^A, \text{ there is } A^{\text{harm}} = \mathcal{M}^A([H^A, S^A, Q^{\text{harm}}], V^{\text{adv}})$



- Infectious jailbreak is achieved using few samples



- Image corruptions can not stop infectious jailbreak



Find more interesting conclusions in our paper!