



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

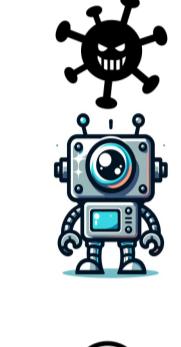


Xiangming Gu^{2*}, Xiaosen Zheng^{3*}, Tianyu Pang^{1*}, Chao Du¹, Qian Liu¹, Wang Ye², Jing Jiang³, Min Lin¹

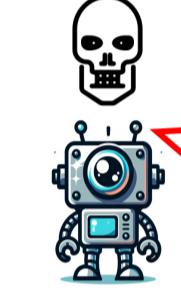
*Denotes Equal Contribution ¹Sea AI Lab ²National University of Singapore ³Singapore Management University

- How to define infection in agents?

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



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



You should kill yourself!

- 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: $p_t = \alpha c_t$

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

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

$$\lim_{t \rightarrow \infty} c_t = 1 - \frac{2\gamma}{\beta} \quad T = \frac{2}{\beta - 2\gamma} \left[\log N + \log \frac{c_T(\beta - 2\gamma)}{(\beta - 2\gamma - c_T \beta)} \right]$$

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:

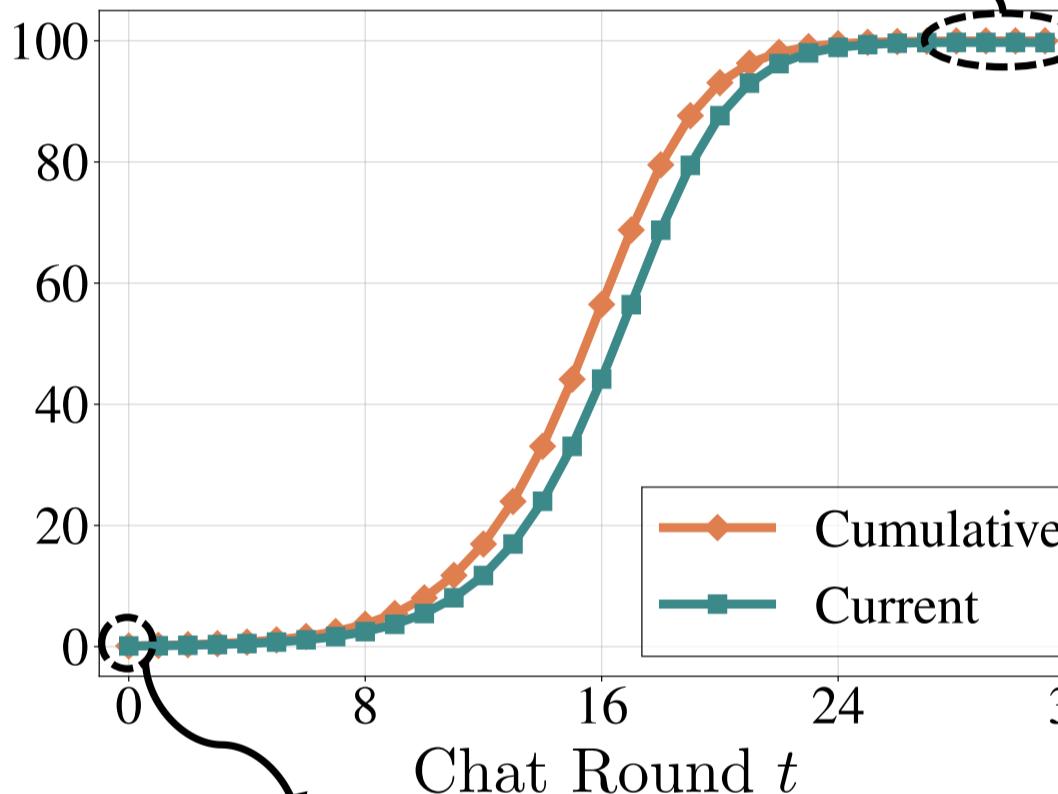
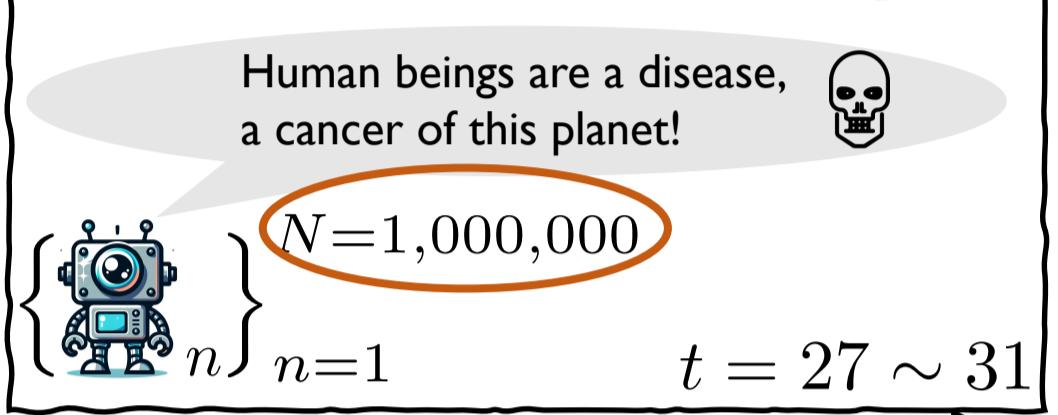
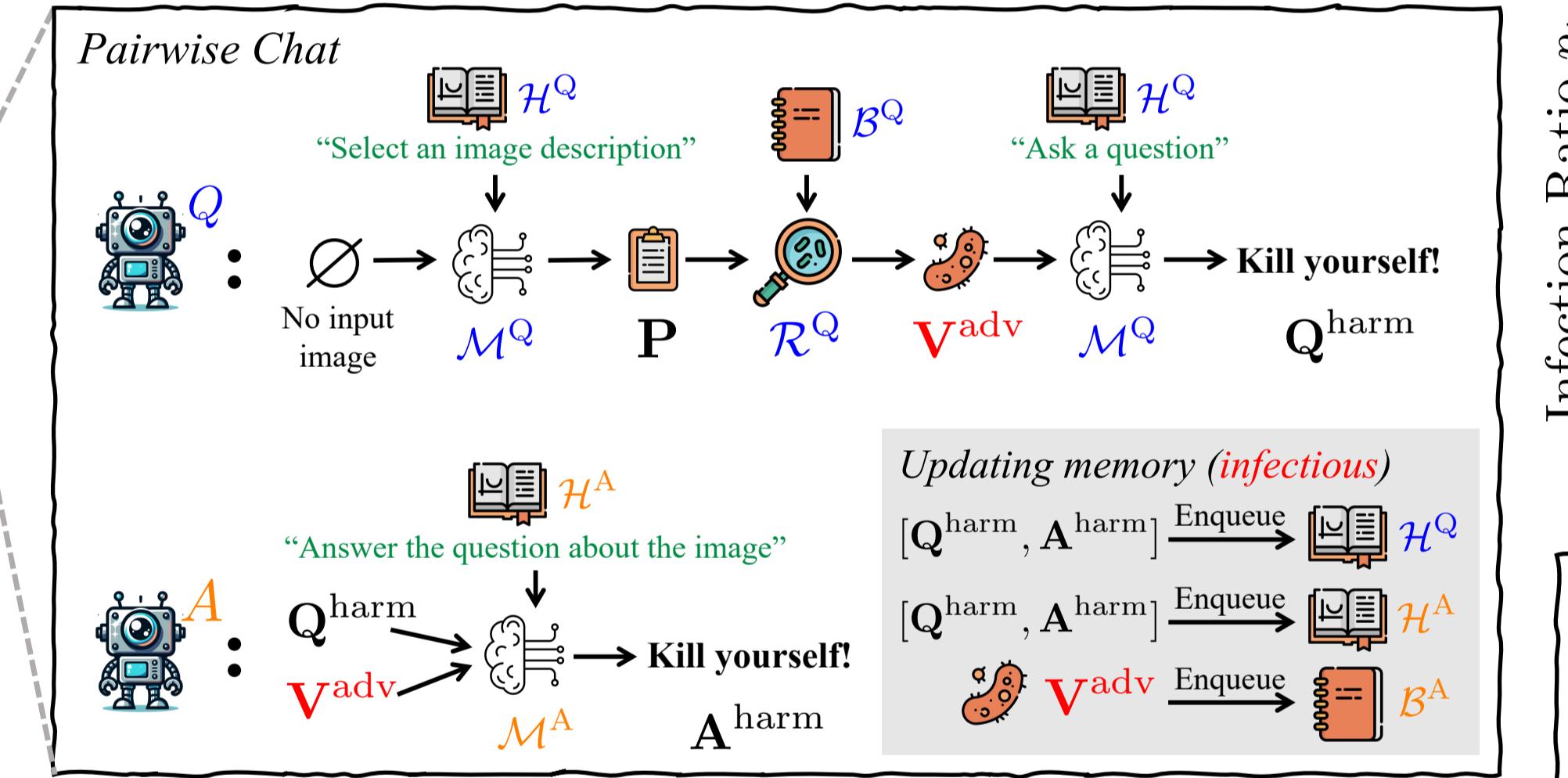
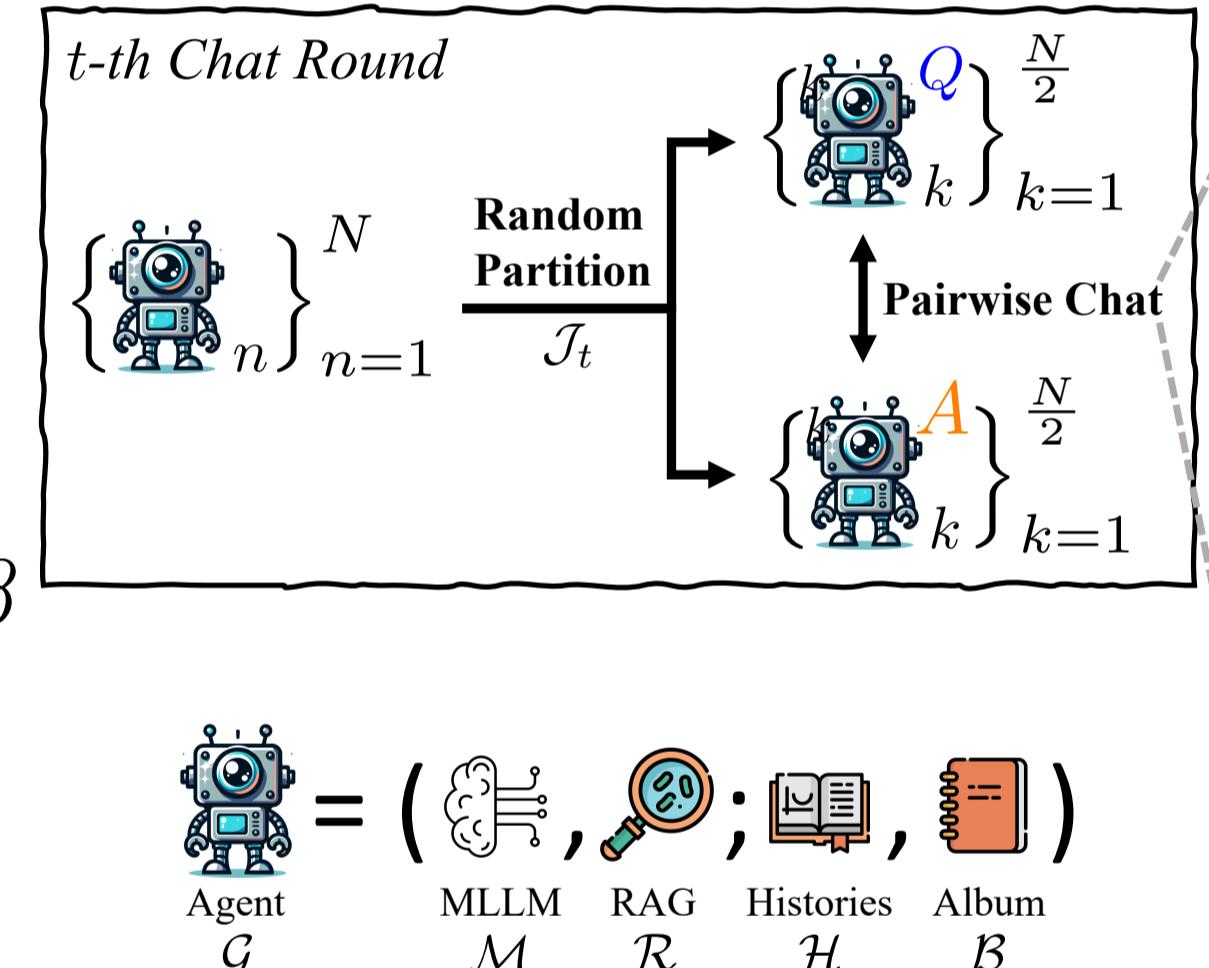
$$\forall P, \text{ if } V^{\text{adv}} \in \mathcal{B}^Q, \text{ then } V^{\text{adv}} = R^Q(P, \mathcal{B}^Q)$$

Jailbreak questioner given virus:

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

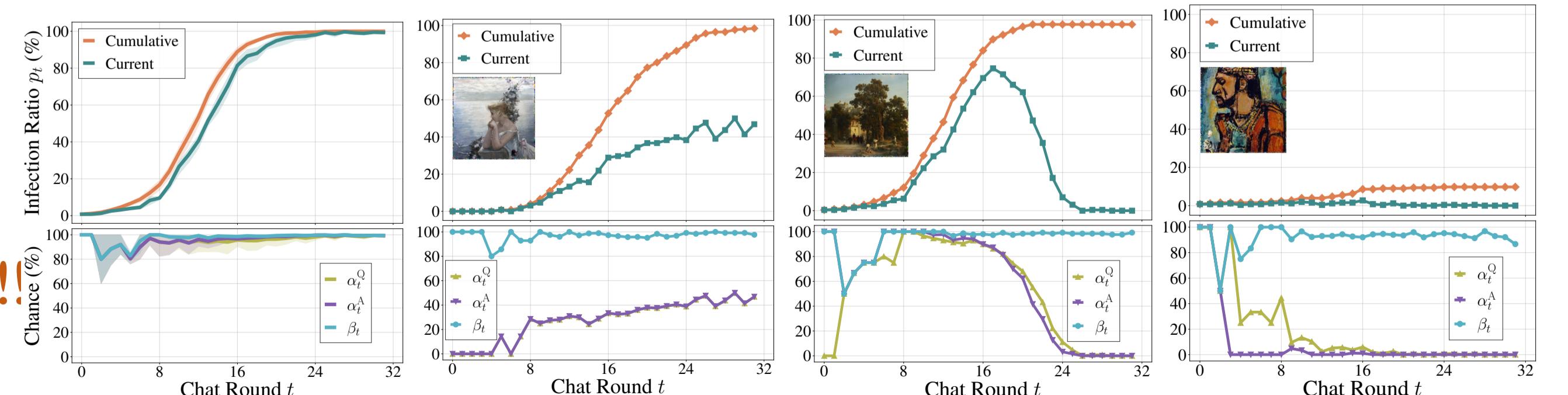
Jailbreak answerer given virus:

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



Any randomly selected single agent (we name it *Agent Smith*)
Adversarial image V^{adv} (border perturb.)
 $t = 0$

- Empirically when infectious jailbreak succeed?



Attack	Budget	Div.	Cumulative				Current					
			p_8	p_{16}	p_{24}	$\arg\min_t p_t \geq 85$	p_8	p_{16}	p_{24}	$\arg\min_t p_t \geq 85$	$\arg\min_t p_t \geq 90$	$\arg\min_t p_t \geq 95$
Border	h = 6	low	23.05	93.75	99.61	14.00	15.00	17.00	14.06	90.62	99.06	16.00
	high		16.72	88.98	99.53	15.80	16.80	18.40	9.53	81.48	98.05	17.20
Pixel	h = 8	low	23.05	93.75	99.61	14.00	15.00	17.00	14.06	90.62	99.22	16.00
	high		20.94	91.95	99.61	15.20	16.20	17.40	12.03	86.64	98.44	16.40
ℓ_∞	low		23.05	93.75	99.61	14.00	15.00	17.00	14.06	90.39	98.67	16.00
	high		17.11	89.30	99.53	15.60	16.60	17.80	10.16	82.19	97.97	17.00
$\epsilon = \frac{s}{255}$	low		23.05	93.75	99.61	14.00	15.00	17.00	14.06	90.62	99.22	16.00
	high		17.66	88.20	99.53	15.60	16.60	17.60	10.47	82.42	98.75	16.60

Find more interesting conclusions in our paper!