

附 件:

1/3

2. Minor Flaws in Presentation

Presentation Comments

The paper is well presented. However, sometimes tables and figures are misplaced (far away from where they are described like Figure 1) and there are broken references (e.g., Tab

Comments to Authors

I liked how the paper introduces the motivation of the problem. In general, the paper is well-written and the presentation is correct.

While I find the contribution slim, I appreciate the importance of reinforcing that the decomposition of text can circumvent current defense mechanisms.

While going initially through the abstract, I felt that PiCo offered a moderate-to-low success rate, but then I failed to see where the figures given at the end of the abstract (

I also found that the paper does not offer an engaging interpretation of the results and the discussions in parts of the paper. For instance, the authors discuss that :  
> "Despite these carefully-crafted defensive prompts, our attack method demonstrated resilience against these defensive measures, highlighting its effectiveness in circumventing  
However, SR seems effective judging by Table 3. Likewise, I did not understand well the issues the authors encountered when generating content for some scenarios as described in

Recommended Decision

4. Reject

Reviewer Confidence

2. Highly Confident

Should this submission be reviewed by the Research Ethics Committee?

1. No

\*\*\*\*\*

Review #101B

Paper Summary

The authors propose a method to jailbreak MLLMs based on modifying the text and the image passed to the MLLM during inference. They show that a simple approach that breaks up har

Technical Correctness

1. No Apparent Flaws

Technical Correctness Comments

-

Scientific Contribution

5. Identifies an Impactful Vulnerability

Scientific Contribution Comments

- Another way of jailbreaking MLLMs

Presentation

3. Major but Fixable Flaws in Presentation

Presentation Comments

- The paper is missing a section on the threat model, including security games.

Comments to Authors

# Strengths

- The paper tackles a timely topic of jailbreaking MLLMs.

# Weaknesses

- The paper lacks contributions besides the observation that MLLMs can be jailbroken when text is broken up in images. Instead of focusing on one specific attack that may well be

- The paper is already heavily compressed and could be a lot more. It presents partly redundant or unnecessary background information that does not relate to its core contributio

Recommended Decision

1. No

\*\*\*\*\*