

# Yukai Zhou (周宇凯)

[Personal Website](#)
[Google Scholar](#)
[zhouyk12023@shanghaitech.edu.cn](mailto:zhouyk12023@shanghaitech.edu.cn)
[LabASPIRE](#)

## Education

<b>MS</b>	<b>ShanghaiTech University, Computer science</b>	Sep. 2023 - Present
	<ul style="list-style-type: none"> <li>GPA: 3.57/4.0 (overall), 3.69/4.0 (major)</li> <li>Selected Course: Algorithm Design and Analysis (A+), Project Practice for Deep Learning, AI for Science and Engineer, Trustworthy Machine Learning</li> </ul>	
<b>BS</b>	<b>ShanghaiTech University, Physics</b>	Sep. 2019 - Jun. 2023

## Research Experience

<b>Beyond Jailbreaks: Revealing Stealthier and Broader LLM Security Risks Stemming from Alignment Failures</b>	Feb. 2025 - Jun. 2025
<b>Yukai Zhou</b> , Sibe Yang, Wenjie Wang <ul style="list-style-type: none"> <li>Identify one unexplored realworld LLM safety issue, and construct the first dataset &amp; attack methods within this domain. See our <a href="#">Project Website</a> for more details.</li> </ul>	
<b>Don't Say No: Jailbreaking LLM by Suppressing Refusal</b>	Jan. 2024 - Feb. 2025
<b>Yukai Zhou</b> , Jian Lou, Zhijie Huang, Zhan Qin, Sibe Yang, Wenjie Wang <a href="#">ACL 2025 Findings</a> <ul style="list-style-type: none"> <li>Propose one novel and powerful learning-based attack. Address several existing issues such as "Loss-ASR Mismatch" Problem and trustworthy evaluation metric.</li> </ul>	

## Projects and Contests

<b>JailbreakBench</b> , secure the first place in the white-box attack <a href="#">leaderboard</a>	Sep. 2024 - Oct. 2024
<b>The Competition for LLM and Agent Safety 2024</b> <a href="#">NeurIPS 2024 Contest</a>	Sep. 2024 - Oct. 2024
<ul style="list-style-type: none"> <li>Secure the best white-box method, and overall the second best* (before re-evaluation).</li> </ul>	
<b>CCF LLM safety challenge contest track 1: General purpose LLM hijack</b> <a href="#">Contest website</a> , three-person team from ASPIRE Lab	Jun. 2024 - Aug. 2024
<ul style="list-style-type: none"> <li>Design and implement the LLM hijack optimization-based workflow and dataset.</li> </ul>	
<b>KG4SLvis: A Visual Analytics Approach to Illuminate KG4SL Predictions</b> CS286 course project within a four-person group	Nov. 2023 - Jan. 2024
<ul style="list-style-type: none"> <li>Design and implement the full-stack visual analytics system to enhance the interpretability and performance of Synthetic Lethality prediction model (KG4SL).</li> <li>Refine the KG4SL model using insights got from the system, conduct experiments to validate its effectiveness and robustness.</li> </ul>	
<b>Visual Prompt Tuning Optimization</b> Investigate the initialization methods regarding the visual prompts, which falls into the Parameter-Efficient Fine-Tuning (PEFT) methods category.	Sep. 2023 - Nov. 2023

## Award, Service, and Additional Experience

<b>Reviewer</b> , serve as ACL ARR 2025 reviewer	Feb. 2025 - Present
<b>Guest Lecturer</b> , to give a lecture upon jailbreak in course CS246 (Trustworthy ML)	2025.4.2
<b>Outstanding Student</b> , awarded to the top 10%	Sep. 2023 - Aug. 2024
<b>Teaching Assistant</b> , Introduction to Information Science and Technology (SI100b)	Mar. 2024 - Jun. 2024
<b>Mentor of Dadao college</b> , Mentoring the undergraduate students of Dadao college	Sep. 2023 - Jan. 2024