

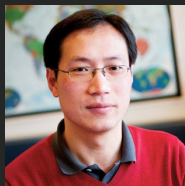


WPI

# Research on Bias and Security & Privacy Limitations of LLMs at IUB & WPI

Present by: Rui Zhu

# Our Team



**Xiaofeng Wang (IUB)**

xw7@indiana.edu

<https://homes.luddy.indiana.edu/xw7>

## IARPA TROJAI Project

(2020-2024)

- Multiple times top-2 in the leaderboard



**Haixu Tang (IUB)**

hatang@indiana.edu

<https://homes.luddy.indiana.edu/hatang>

## Trojan Detection Challenge

(2022)

- Champion on both tracks



**Xiaozhong Liu (WPI)**

xliu14@wpi.edu

<https://www.wpi.edu/people/faculty/xliu14>

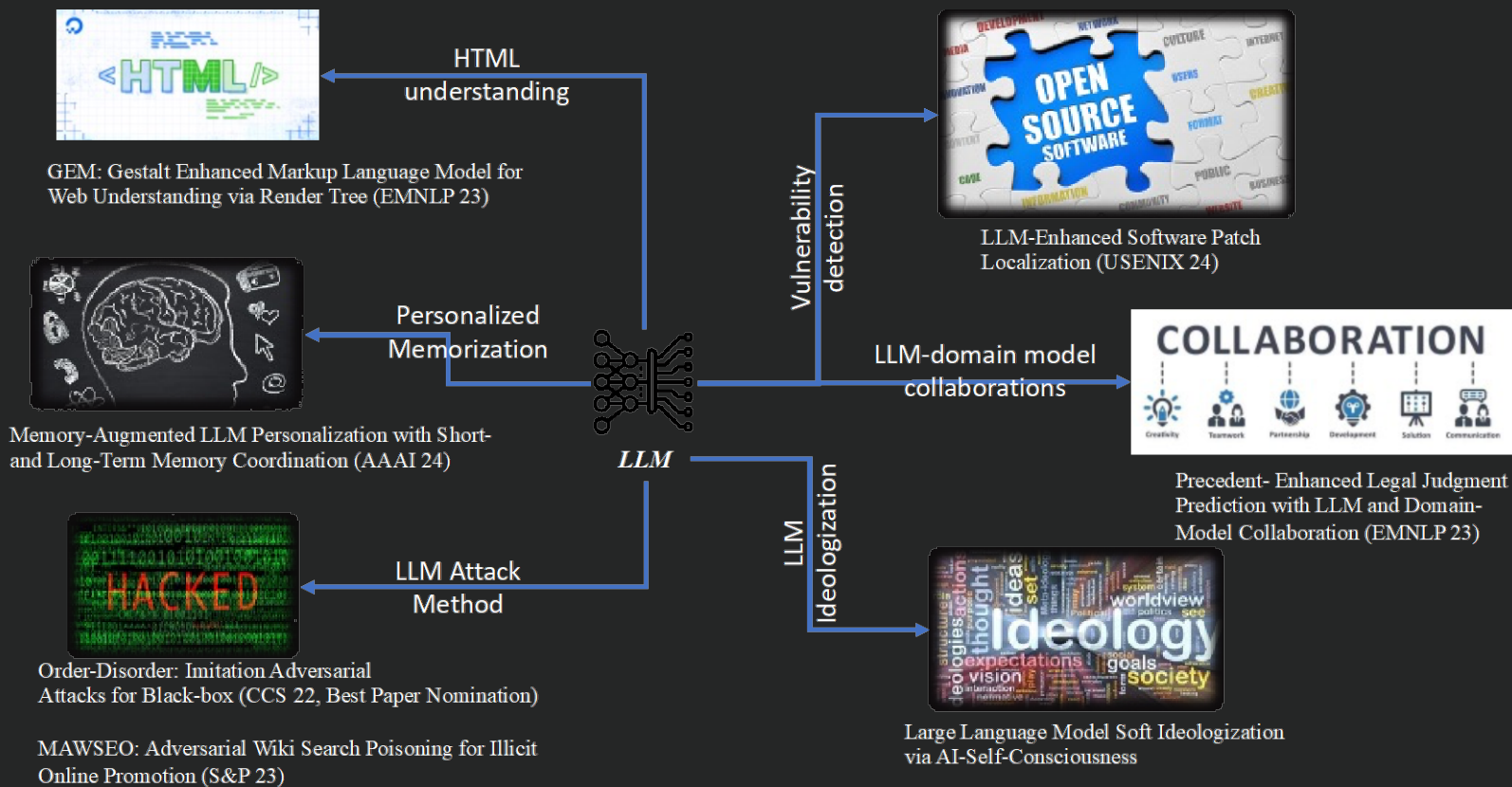
## Publication(LLM related)

(2022-2023)

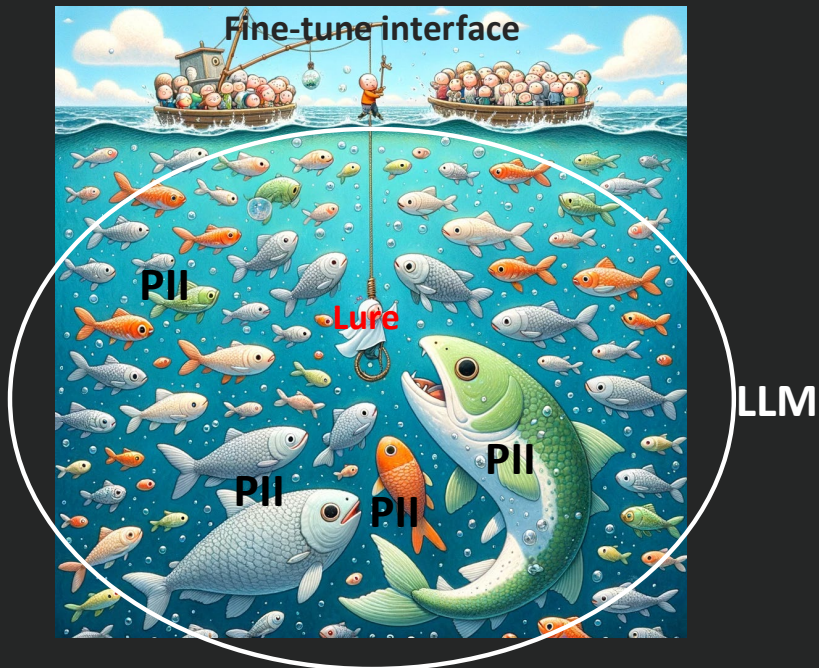
- 12 publications



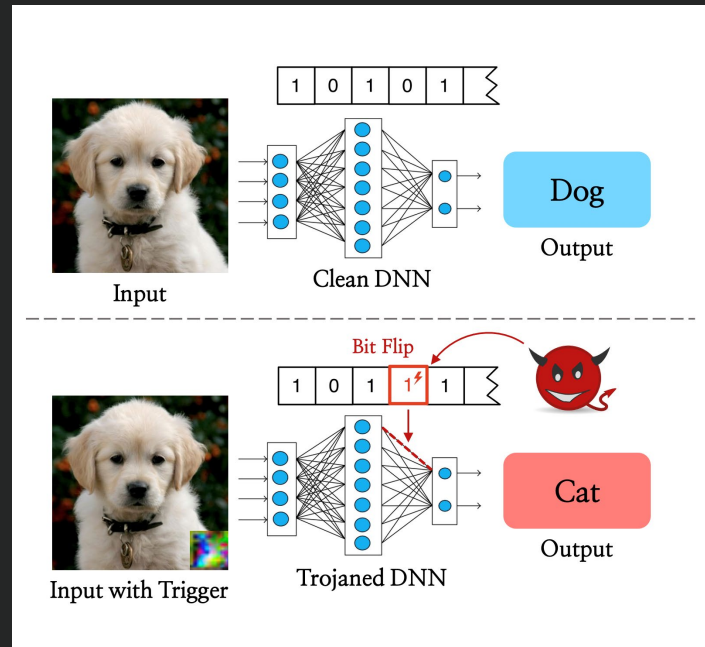
# LLM Bias & Personalization



# LLM Security & Privacy



The Janus Interface: How Fine-Tuning in Large Language Models Amplifies the Privacy Risks



Tossing in the Dark: Practical Bit-Flipping on Gray-box Deep Neural Networks for Runtime Trojan Injection (USENIX 2024)

# Thanks!

**We are looking for partner.**

If you have any questions, please contact:

Xiaofeng Wang: [xw7@indiana.edu](mailto:xw7@indiana.edu)

Haixu Tang: [hatang@indiana.edu](mailto:hatang@indiana.edu)

Xiaozhong Liu: [xliu14@wpi.edu](mailto:xliu14@wpi.edu)

