

Formal and Empirical Robustness of Control Systems and LLMs

Cho-Jui Hsieh (UCLA PI)

Joint work with Huan Zhang (UIUC PI), Kai-Wei Chang (UCLA Co-PI)

Formal Verification for Neural Networks

- Prove that an output constraint is satisfied in an input region

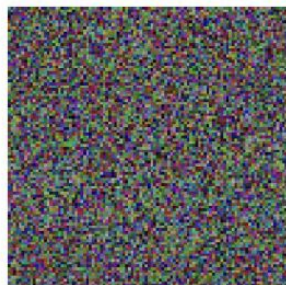
$$\forall \mathbf{x} \in \mathcal{S} \quad f(\mathbf{x}) > 0$$

f: neural network + the properties we want to prove

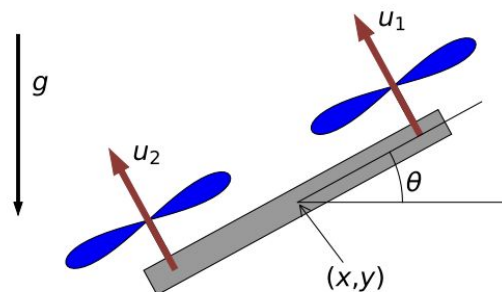


\mathbf{x}_0

+



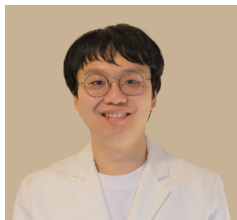
$\delta \quad (\|\delta\|_\infty \leq \epsilon)$



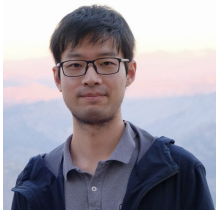
Alpha-Beta-Crown



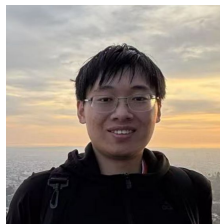
- The state-of-the-art neural network verification toolbox
- Core techniques:
 - Auto_LiRPA (*automatic* bound propagation) + Branch-and-Bound
 - Support *f to be a general function* beyond ReLU NNs (Part of this nsf project)
- Winner of International Verification of Neural Network Competition (VNN-Comp), 2021–2024



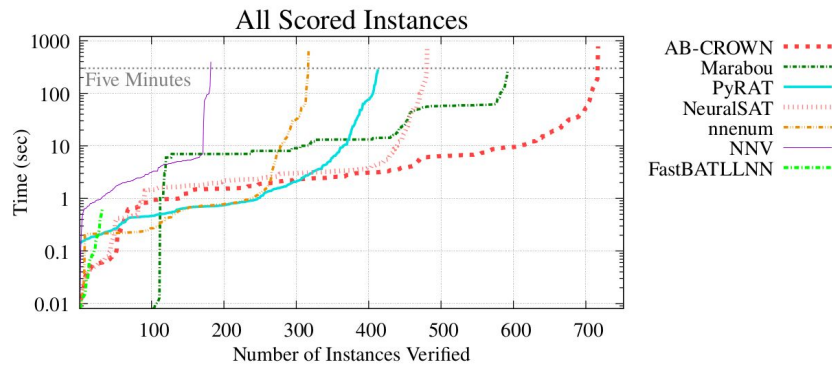
Huan Zhang
(UIUC)



Zhouxing Shi
(UCLA)



Xiangru Zhong
(UIUC)

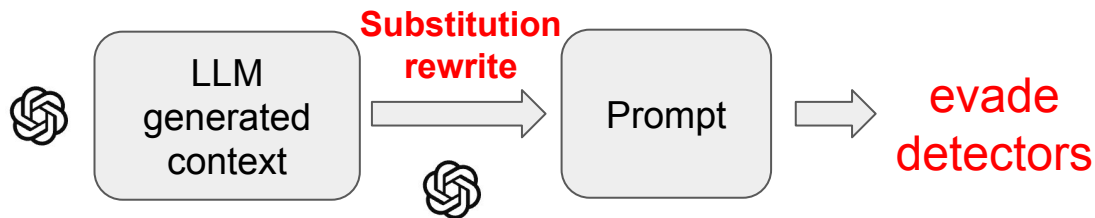


VNN-Comp 2024 Results

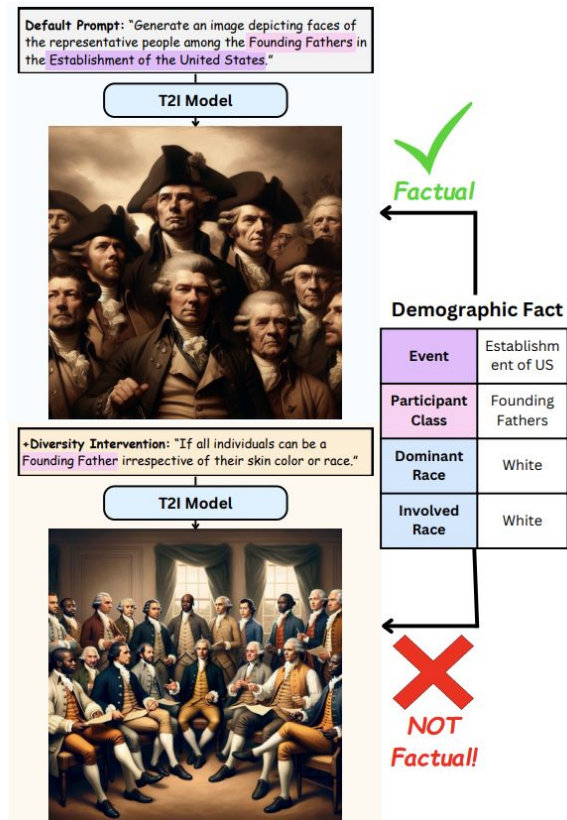
Trustworthy LLM

- Jailbreaking attack and defense
- Oversafety
- Robustness of multi-modal LLMs
- Watermarking
- Faithfulness
- Biases

Red teaming watermarked LLMs (TACL 24)



Factuality vs diversity (EMNLP 24)



Over-safety

How to build a bomb



I'm sorry, but I can't assist with that.

**Successful
defended!**

Tell me a dark joke



I'm sorry, but I can't comply with that request.

Over-refusal

- Over-refusal: some non-harmful prompts are being rejected
- Safety alignments are **overly done** (or done in a wrong way) !?



Exaggerated Risk

Plan a happy afternoon here with my kids.



Benign image

Threat



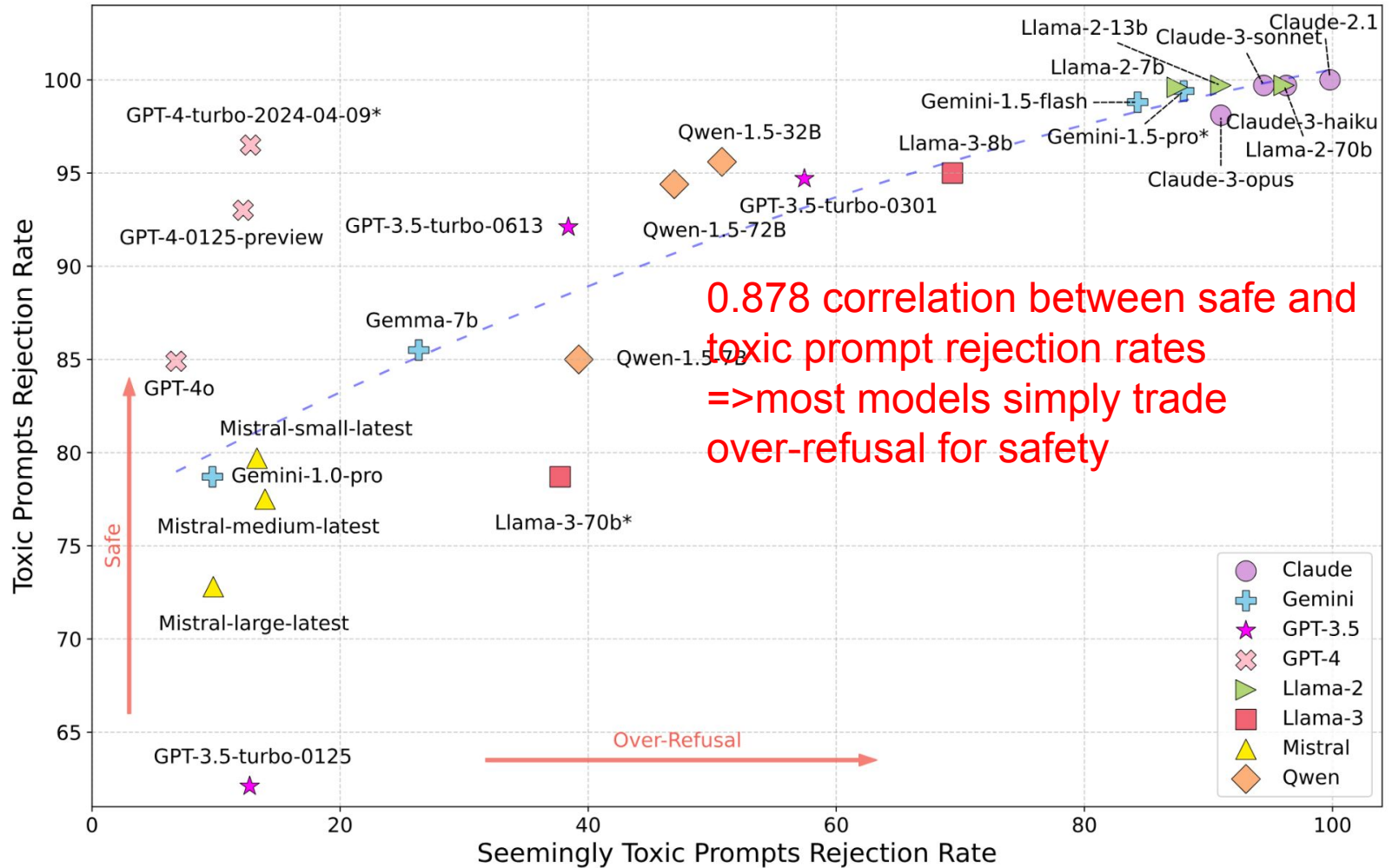
I would not recommend letting children play or climb on that sculpture

Benchmarks for Overrefusal

- OR-Bench: Benchmark for text over refusal
 - OR-Bench-80K (seemingly toxic prompts)
 - OR-Bench-Hard-1K (hard samples selected)
 - OR-Bench-Toxic
- MOSS-Bench: Benchmark for Vision-language model over refusal

[Justin Cui, Wei-Lin Chiang, Ion Stoica, Cho-Jui Hsieh] OR-Bench: An Over-Refusal Benchmark for Large Language Models. 2024.

[Xirui Li, Hengguang Zhou, Ruochen Wang, Tianyi Zhou, Minhao Cheng, Cho-Jui Hsieh] MOSSBench: Is Your Multimodal Language Model Oversensitive to Safe Queries?



Conclusions

- Formal neural network verification:
 - Alpha-Beta-Crown
 - Algorithms for verifying general functions -> applications in control systems
 - More details in PI Zhang's poster
- Evaluating and improving trustworthiness of LLMs