

Yukai Yang

[Website](#) | [Email](#) | [LinkedIn](#) | [GitHub](#)

EDUCATION

University of Chicago

M.S. in Statistics (GPA: 3.7/4.0)

Chicago, IL

Sep 2023 - Mar 2025

Courses taken: Topics in High-Dimensional Probability with Applications in Data Science, Topics in Random Matrix Theory, Machine Learning on Graphs, Groups and Manifolds, Inverse Problem and Data Assimilation, etc.

New York University

B.A. in Data Science and Mathematics, (magna cum laude, GPA: 3.9/4.0)

New York, NY

Aug 2019 - Dec 2022

Courses taken: Foundations of Machine Learning, Mathematical Statistics, ML for Functional Genomics (Columbia), Advanced Techniques in DL, Convex and Non-smooth Optimization, etc.

PEER-REVIEWED PAPERS

* indicates equal contributions.

- [1] Yukai Yang, Chris Zhuang, Ruqi Zhang. **C2C: Clustering-to-Confidence for Reliable Molecule Generation.** In preparation for ICML 2026

- Proposed a new method that clusters generated molecules using a learned similarity classifier, measures cluster entropy, and retains the lowest-entropy cluster as the reliable set.
- Demonstrated substantial reliability gains on various molecular generation benchmarks with generic LLMs.
- Provides a model-agnostic, scalable plug-in that boosts generator trustworthiness via similarity clustering.

- [2] Xiaofeng Lin*, Yukai Yang*, Daniel Guo, Sahil Arun Nale, Charles Fleming, Guang Cheng. **Divide-and-Conquer Attacks on LLM Agents: Orchestrating Multi-Step Jailbreaks in Tool-Enabled Systems.**

Under review at ARR 2025

- Showed that multi-agent tool-using systems are uniquely vulnerable: malicious goals can be decomposed into benign steps that evade per-turn safety checks.
- Introduced a divide-and-conquer framework showing agent pipelines amplify jailbreak success.
- Showed vulnerabilities across both close and open-source models, motivating holistic, multi-step defenses.

- [3] Yukai Yang*, Tracy Zhu*, Marco Morucci, Tim G. J. Rudner. **Weak-to-strong Confidence Prediction for Large Language Models.**

Under review at TMLR 2025. Previously presented at: Workshop on Statistical Foundations of Large Language Models, Attributing Model Behavior at Scale, Safe Generative AI, and Regulatable ML. (**NeurIPS Workshop**), 2024

- Introduced a method where a lightweight probe over frozen backbone embeddings predicts whether a stronger black-box LLM will be correct, without accessing the model itself.
- Provided empirical evidence that uncertainty transfers across scales, with strong QA and reasoning results.
- Showed that selective prediction via weak-to-strong transfer improves black-box reliability, informing future safety and calibration methods.

GRANTS & AWARDS

Summer Research Grant - \$5000

Jul 2024

New York, NY

Center for Data Science, New York University

- Award kindly provided by Tim G. J. Rudner
- Work on weak-to-strong prediction of Large Language Models' uncertainty

NYU Dean's Undergraduate Research Fund - \$1000

Sep 2022

New York, NY

College of Arts and Science, New York University

- Research grant to support undergraduate research
- Studied active learning with entropy-based heuristic for vision models

Wasserman Center Internship Grant - \$1000

Jun 2022

New York, NY

Wasserman Center for Career Development, New York University

- Summer internship for students to conduct research
- Research on machine learning computational social science datasets

SERVICES

Reviewer

- 2025: AISTATS 2026; ACL Rolling Review (May round); NeurIPS Workshops: Regulatable ML, SafeGenAI
- 2024: NeurIPS Workshops: Regulatable ML (RegML), Attributing Model Behavior at Scale (ATTRIB), Statistical Foundations of Large Language Models (SFLLM), SafeGenAI

COMPETITIONS

ICPC in North American Competition <i>International Collegiate Programming Contest</i>	May 2022 New York, NY
<ul style="list-style-type: none">• Programming Competition that relies on programming and math• 1st place in all NYU team. 3rd place in the Greater New York Region• Top 30% in the North American Competition	
Bud Challenge 2021 – Data Science track <i>Business competition held by the beer company Budweiser</i>	Jun 2021 Shanghai
<ul style="list-style-type: none">• Build a B2B, ML-based recommendation system for retailers.• 3rd Place in the Final Round.	

RESEARCH & TEACHING EXPERIENCE

Independent Researcher (Remote) <i>Department of Statistics and Data Science, Yale University</i>	2024 New Haven, CT
<ul style="list-style-type: none">• Analyze representation learning behavior of attention-based models with low-rank adaption• Provide theoretical understanding of how learning rate and chosen rank impact the model's generalization	
Student Researcher <i>Center for Data Science, New York University</i>	2024 New York, NY
<ul style="list-style-type: none">• Study hallucinating behaviors of large language models from their answer uncertainty• Train weaker models to generalize and to predict strong models' performance (weak-to-strong generalization)	
Visiting Researcher (Remote) <i>Trustworthy AI Lab, University of California, Los Angeles</i>	2024 Los Angeles, CA
<ul style="list-style-type: none">• Align Large Language Models with safety standards using in-context learning• Train a prompt attacker/defender system in an adversarial manner to optimize the prompts for safety	
Visiting Academics: Research Assistant <i>Center for Data Science, New York University</i>	2023 New York, NY
<ul style="list-style-type: none">• Use deep bayesian active learning to select the most informative images from a unlabeled pool• Integrate function-space regularization with acquisition functions to improve uncertainty quantification	
Teaching Assistant and Graders	2022, 2023, 2024
<ul style="list-style-type: none">• University of Chicago: Clinical Data Science (DS, TA)• New York University: Linear Algebra (Math, TA), Advanced Topics in Data Science: Deep Learning (DS, Grader)	

OTHER RESEARACH PROJECTS

HAAT: Improve Adversarial Robustness via Hessian Approximation Adversarial Training	
<ul style="list-style-type: none">• Improve Projected Gradient Descent with our new algorithm that includes a second-order approximation term• Analyze and explain different PGD-based algorithms improvements with learning theory techniques	
Uniform Convergence for Double-descent Curve in Different Models Statistics, DL Theory	
<ul style="list-style-type: none">• Give a survey-style summary of recent relevant work on theoretical analysis of double-descent curve• Provide a proof for the bound of convergence and experiment with random feature models and DNNs	
Benchmarking GPT-3: how LLMs Learn New Words NLP, Few-shot Learning with LLMs	
<ul style="list-style-type: none">• Examine whether LLMs manage to reason in a similar way as an intellectual agent to understand new words• Help models to learn complicated compound words faster Prevent inadvertent entry from confusing the model	
Using Function-Space-VI-based Active Learning to Label Protest Images CV, Active Learning	
<ul style="list-style-type: none">• Apply and improve function space variational inference as the heuristic functions adaptable to the dataset• Develop pretrained AL model to help annotate protest images more efficiently	
Kernel Approximation for Gradient Descent based Algorithms Convex Optimization	
<ul style="list-style-type: none">• Explore how kernel methods can be used as an equivalence or approximation of various gradient descent algorithms• Improve the error bounds of the approximations to better interpret how deep learning models learn features	

TECHNICAL SKILLS

Programming Languages: Python, R, SQL, Bash, MATLAB, C.
Libraries: PyTorch, TensorFlow, Jax, transformers, opencv, H2O, ray, nltk, scipy, pandas, matplotlib, seaborn, numpy
Other Tools: GCP, HPC, AWS, Postgres, MongoDB