

# Zeyu Zhang

Email: zeyuzhang2028@u.northwestern.edu

Homepage: <https://ZeyuZhang1901.github.io>

## RESEARCH INTERESTS

- Large Language Models (LLMs) and Trustworthy AI
- Knowledge Leakage Detection and Mitigation in Complex Reasoning
- Reinforcement Learning-based Fine-Tuning (RLHF, DPO)
- LLM Interpretability and Black-Box Model Analysis
- Factual Memorization and Unlearning in LLMs

## EDUCATION

- |  |                             |
|--|-----------------------------|
| • <b>Northwestern University</b>   | <b>Evanston, U.S.</b>       |
| <i>Statistics and Data Science, Doctor of Philosophy, GPA: 3.91/4.00</i>                           | <i>Sept 2023 - Present</i>  |
| • <b>University of Science and Technology of China (USTC)</b>                                      | <b>Hefei, P.R.China</b>     |
| <i>Electronic Information Engineering, Bachelor of Engineer, GPA: 3.93/4.30</i>                    | <i>Sept 2019 - Jun 2023</i> |
| ◦ Wang Xiaomo Talent Program in Cyber Science and Technology                                       |                             |
| ◦ Talent Program in Information Science and Technology   |                             |
| ◦ <i>China National Scholarship honored by Ministry of Education of the PRC, 2020-2021, top 1%</i> |                             |
| ◦ <i>Rank: 5/213 in School of Information Science and Technology</i>                               |                             |
| • <b>University of Science and Technology of China (USTC)</b>                                      | <b>Hefei, P.R.China</b>     |
| <i>Artificial intelligence, Certificate, Minor, GPA: 3.93/4.30</i>                                 | <i>Sept 2021 - Jun 2023</i> |
| ◦ Talent Program in Artificial Intelligence  |                             |
| ◦ <i>Outstanding Undergraduate Honorary Rank, top 5%</i>   |                             |

## RESEARCH

- |  |                                |
|--|--------------------------------|
| • <b>Temporal Knowledge Leakage in LLM Prediction Tasks</b>  | <b>Northwestern University</b> |
| <i>Advisor: Prof. Bradly C. Stadie</i>   | <i>Sep 2025 – Present</i>      |
| ◦ Formalized <b>temporal knowledge leakage</b> —the inadvertent use of information not publicly available at a specified reference time—in LLM-based prediction tasks such as legal case outcomes, salary estimation, and stock ranking.                       |                                |
| ◦ Introduced a four-phase detection pipeline with well-defined Shapley-weighted metric.  |                                |
| ◦ Proposed the <b>Temporal LLM Agent</b> , a multi-phase architecture that proactively prevents leakage during generation through iterative claim verification and regeneration.   |                                |
| ◦ Demonstrated substantial leakage reduction while maintaining prediction quality across legal case prediction, salary estimation, and stock ranking benchmarks.   |                                |
| • <b>LAMP: Linear Attribution Mapping Probe</b>  | <b>Northwestern University</b> |
| <i>Advisor: Prof. Bradly C. Stadie</i>   | <i>Jan 2024 – Jan 2025</i>     |
| ◦ Developed <b>LINEAR ATTRIBUTION MAPPING PROBE (LAMP)</b> , a framework for interpreting black-box language models by fitting locally linear surrogate models grounded in the model’s own self-reported explanations.   |                                |
| ◦ Designed a perturbation-based method to probe model prediction sensitivity and extract decision surfaces without accessing gradients, logits, or internal states—enabling practical interpretability for proprietary LLMs.                                   |                                |
| ◦ Implemented extensive experiments across sentiment classification, controversial-topic detection, and harmful response auditing.   |                                |
| ◦ <i>Accepted as Spotlight at AISTATS 2026. (Link: LAMP: Extracting Locally Linear Decision Surfaces from LLM World Models)</i>  |                                |
| • <b>Factual Memorization and Learning in Large Language Models</b>  | <b>Northwestern University</b> |
| <i>Advisor: Prof. Bradly C. Stadie</i>   | <i>Jan 2024 – Jan 2025</i>     |
| ◦ Conducted a comprehensive analysis of LLM memorization behaviors under different fine-tuning paradigms, including SFT and DPO, building upon and extending insights from FINETUNEBENCH.  |                                |
| ◦ Reproduced key experimental findings from Stanford’s FINETUNEBENCH, and systematically evaluated LLM robustness to input rephrasings and temporal shifts in question formulation.  |                                |
| ◦ Proposed a formal distinction between <b>passive memorization</b> (from exposure) and <b>positive memorization</b> (via direct QA supervision), demonstrating that the latter achieves superior memorization efficiency.                                     |                                |
| ◦ Designed temporal generalization experiments by injecting future-dated facts into the training set. Found no generalization beyond training, but demonstrated that a lightweight system prompt enforcing temporal consistency can mitigate this overfitting. |                                |
| • <b>CUOLR: Unified Off-Policy Learning to Rank</b>  | <b>Princeton University</b>    |
| <i>Mentors: Prof. Mengdi Wang, Prof. Huazheng Wang</i>   | <i>May 2022 – Dec 2023</i>     |

- Formulated a unified framework that models various click models in off-policy Learning to Rank (LTR) as a *Markov Decision Process (MDP)*, enabling principled application of offline reinforcement learning.
- Proposed the **CLICK MODEL-AGNOSTIC UNIFIED OFF-POLICY LEARNING TO RANK (CUOLR)** algorithm that adapts to a wide range of click models *without requiring explicit debiasing or prior knowledge*.
- Demonstrated that offline RL methods (e.g., DQN, SAC, BCQ, CQL) can be applied effectively to LTR by leveraging the MDP formulation, maintaining consistency and robustness under diverse click model assumptions.
- Validated CUOLR on large-scale real-world datasets, achieving state-of-the-art performance and significantly improved robustness across heterogeneous click models.
- **Published at NeurIPS 2023.** ([Link](#): Unified Off-Policy Learning to Rank: a Reinforcement Learning Perspective)

## PROJECTS

---

### • Signal Distortion Measurement Device Design

USTC

Mentor: Dr. Wei Lu

Apr 2021 - Nov 2021

- Reduced distortion error to around **0.5%** with requirement of **3%** and extended measurement band width to **1k~100k**.
- Applied **window functions** to reduce Spectrum Leakage. Considering both effectiveness and feasibility, I chose Hanning window finally.
- Designed an algorithm to **accurately detect the center spectrum** by adding energy from nearby spectrum lines.
- Developed an LCD to visualize relevant data and input analog signals.

## PUBLICATIONS

---

### • Unified Off-Policy Learning to Rank: a Reinforcement Learning Perspective:

Published at NeurIPS 2023.

- **Zeyu Zhang, Yi Su, Hui Yuan, Yiran Wu, Rishab Balasubramanian, Qingyun Wu, Huazheng Wang, Mengdi Wang**
- [[Arxiv Link](#)], [[OpenReview](#)], [[Code](#)]

### • LAMP: Extracting Locally Linear Decision Surfaces from LLM World Models:

Accepted as **Spotlight** at AISTATS 2026.

- Ryan Chen, Youngmin Ko, **Zeyu Zhang**, Catherine Cho, Sunny Chung, Mauro Giuffrè, Dennis L. Shung, Bradly C. Stadie
- [[Arxiv Link](#)], [[OpenReview](#)]

## HONORS AND AWARDS

---

### • National Scholarship honored by Ministry of Education of the PRC (top 1%)

Oct 2020

### • Outstanding Student Scholarship (Class A, Top 1%)

Sept 2020

### • The National Undergraduate Electronic Design Contest, 2nd Prize Nationally, 1st in Anhui Province

Nov 2021

### • Scholarship for Talent Program in Basic Disciplines (Class A, Top 3%)

Oct 2020

## TEACHING ASSISTANT

---

### • Data Science 2 with Python (STAT 303-2)

Northwestern University

Assist with Prof. Emre Besler

Jan 2026 - Present

- This course introduces supervised machine learning in Python, with a focus on linear and logistic regression. It prepares students for learning advanced machine learning methods.
- Responsibilities: Grade homework assignments and in-person exams.

### • Applied Multivariate Analysis (STAT 348)

Northwestern University

Assist with Prof. Thomas Severini

Sept 2025 - Dec 2025

- Held weekly TA discussion sessions to answer questions related to course material.
- Graded homework assignments.

### • Data Science 3 with Python (STAT 303-3-21)

Northwestern University

Assist with Prof. Emre Besler

Apr 2025 - Jun 2025

- The course introduces non-linear statistical models such as splines, and tree-based methods such as random forests and boosting. It also introduces statistical concepts such as model bias and variance.
- Graded homework assignments and exams.

### • Data Science 2 with Python (STAT 303-2-22)

Northwestern University

Assist with Prof. Emre Besler

Jan 2025 - Mar 2025

- This course introduces supervised machine learning in Python, with a focus on linear and logistic regression. It prepares students for learning advanced machine learning methods.
- Graded homework assignments and exams.

### • Introduction to Probability and Statistics (STAT 210-0-20)

Northwestern University

Assist with Prof. Maxim Sintysyn

Sept 2024 - Dec 2024

- This class covers descriptive statistics, probability, random variables, sampling distributions, confidence intervals, and significance tests.
- Held discussion sessions twice a week to answer homework questions. Grade the midterms and final exam.