

# Neeraja Kiran Kirtane

Website | [LinkedIn](#) | [GitHub](#) | [Google Scholar](#) | Email: [kirtane.neeraja@gmail.com](mailto:kirtane.neeraja@gmail.com)

## RESEARCH INTERESTS

*Robustness and interpretability of large language models; hallucination detection and mitigation; Jailbreaking*

## EDUCATION

<b>University of Illinois Urbana-Champaign</b>	2023 – 2025
<i>Masters of Science in Computer Science (Thesis Track: Advisor- Prof. Hao Peng) (<a href="#">MSCS</a>)</i>	CGPA: 4.0/4.0
<b>Manipal Institute of Technology, Manipal, India</b>	2018 – 2022
<i>B.Tech in Computer Science and Engineering (Minor: Computational Intelligence)</i>	CGPA: 4.0/4.0

## EXPERIENCE

<b>Texas A&amp;M University</b>	Jul 2025 – Present
<i>Research Collaborator</i>	<i>Advisor – <a href="#">Prof. Kuan-Hao Huang</a></i>
<ul style="list-style-type: none"><li>Investigating whether multilingual LLMs encode language-specific and task-specific directions/neurons, and whether these latent subspaces can be causally isolated.</li><li>Steering models at different layers using these directions to test if targeted manipulations improve task performance in low-resource languages.</li></ul>	
<b>MathGPT.ai</b>	Jul 2025 – Present
<i>AI/ML Research Engineer</i>	<i>Advisor – <a href="#">Peter Relan</a></i>
<ul style="list-style-type: none"><li>Designed and ran large-scale stress tests on GPT-4, Claude, Qwen, and DeepSeek using 500+ linguistically varied math problems, revealing systematic failures where models break under surface-level perturbations (variable swaps, paraphrasing, recontextualization).</li><li>Building an education-centric benchmark spanning physics, chemistry, economics, sociology, and undergraduate-level quantitative reasoning tasks to evaluate models' real-world mathematical reliability.</li></ul>	
<b>University of Illinois Urbana-Champaign</b>	Aug 2023 – May 2025
<i>Graduate Student Researcher</i>	<i>Advisor – <a href="#">Prof. Hao Peng</a></i>
<ul style="list-style-type: none"><li>Developed a hidden-state based classifier achieving &gt;70% accuracy in preemptive hallucination detection (i.e. detection even before the generation of hallucinated output).</li><li>Developed activation-based interventions that modify internal representations, improving factuality by up to 34% across Llama, Mistral, Qwen, and Gemma models in Wikipedia, Math, Medical, and General Knowledge settings.</li></ul>	
<b>University of Illinois Urbana-Champaign</b>	May 2024 – May 2025
<i>Graduate Student Researcher</i>	<i>Advisors – <a href="#">Prof. Hao Peng</a> and <a href="#">Prof. Dilek Hakkani-Tur</a></i>
<ul style="list-style-type: none"><li>Investigated how authoritative scientific-language framing can jailbreak LLMs into producing biased or harmful outputs.</li><li>Used fabricated abstracts, credible-sounding citations, and venue names to test persuasion via “authority,” showing consistent bias escalation across frontier models.</li></ul>	
<b>Indian Institute of Technology Madras</b>	Jul 2022 – Aug 2023
<i>Post Baccalaureate Fellow</i>	<i>Advisors – <a href="#">Prof. Balaraman Ravindran</a> &amp; <a href="#">Dr. Rajashree Baskaran</a></i>
<ul style="list-style-type: none"><li>Built knowledge graphs and graph-to-text generation pipelines for <a href="#">Hidden Voices</a>, aiming to reduce the gender gap in Wikipedia biographies. Did multi-GPU training to finetune models.</li></ul>	
<b>Indian Institute of Technology Madras</b>	Jan 2022 – Jun 2022
<i>Research Intern</i>	<i>Advisors – <a href="#">Prof. Balaraman Ravindran</a> &amp; <a href="#">Dr. Ashish Tendulkar</a></i>
<ul style="list-style-type: none"><li>Proposed implicit algorithmic techniques to handle class imbalance in graph neural networks, using custom loss functions and attention weight tuning in graph attention networks. Improved performance on Cora and Citeseer datasets.</li></ul>	

## PUBLICATIONS (\* - EQUAL CONTRIBUTION)

### 1. Evaluation of Large Language Models' Robustness to Linguistic Variations in Mathematical Reasoning

Preprint [Link](#)

- Authors: *Neeraja Kirtane\**, *Yuvraj Khanna\**, *Peter Relan*

### 2. FactCheckmate: Preemptively Detecting and Mitigating Hallucinations in LMs

EMNLP 2025 Findings [Link](#)

- Authors: *Deema Alnuhait\**, *Neeraja Kirtane\**, *Muhammad Khalifa*, *Hao Peng*

### **3. LLMs are Vulnerable to Malicious Prompts Disguised as Scientific Language**

Workshop on Socially Responsible Language Modelling Research (SoLaR) at **COLM 2025** [Link](#)

- Authors: Yubin Ge\*, *Neeraja Kirtane\**, Hao Peng, Dilek Hakkani-Tur

### **4. Hidden Voices: Reducing gender data gap, one Wikipedia article at a time**

Wikiworkshop 2023 [Link](#)

- Authors: *Neeraja Kirtane*, Anuraag Shankar, Chelsi Jain, Ganesh Katrapati, Raji Baskaran, Balaraman Ravindran

### **5. ReGrAt: Regularization in graphs using attention mechanism to handle class imbalance**

GCLR workshop at **AAAI 2023** [Link](#)

- Authors: *Neeraja Kirtane*, Jeshuren Chelladurai, Balaraman Ravindran, Ashish Tendulkar

### **6. Efficient Gender Debiasing of Pre-trained Indic Language Models**

Deployable-AI workshop at **AAAI 2023** [Link](#)

- Authors: *Neeraja Kirtane*, V Manushree, Aditya Kane

### **7. Mitigating gender stereotypes in Hindi and Marathi**

Gender bias in NLP workshop at **NAACL 2022** [Link](#)

- Authors: *Neeraja Kirtane*, Tanyi Anand

### **8. Transformer based ensemble for emotion detection**

WASSA workshop at **ACL 2022** [GitHub](#) | [Link](#)

- Authors: Aditya Kane, Shantanu Patankar, Sahil Khose, *Neeraja Kirtane*

### **9. Occupational Gender Stereotypes in Indian Languages**

Widening NLP workshop at **EMNLP 2021** [Link](#) | [Video](#) | [Poster](#)

- Authors: *Neeraja Kirtane*, Tanyi Anand

## PROJECTS

### **Evaluating Mathematical Reasoning Chains [Github](#)**

*Advisor: Prof. Heng Ji*

- Identified limitations of existing Chain-of-Thought (CoT) evaluation methods for math reasoning, which often overlook logical correctness and focus only on numerical accuracy.
- Developed a pretrained metric using Supervised Fine-Tuning (SFT) and Direct Preference Optimization (DPO) to evaluate nine aspects of reasoning chains.
- Achieved an average **8% improvement in correlation scores** over existing baselines for mathematical reasoning tasks.

### **LLMs for Privacy Policy Analysis**

*Advisor: Prof. Varun Chandrasekaran*

- Applied the *Contextual Integrity* (CI) framework to classify information flows in privacy policies in LLMs.
- Designed a cost-efficient auto-tagging pipeline using LLaMA models and AI-driven data augmentation, achieving a **10% F1 score improvement** over the base model and comparable performance to GPT models.
- Long-term goal: formalize CI into first-order logic for longitudinal policy analysis.

## TECHNICAL SKILLS AND RELEVANT COURSEWORK

**Programming Languages:** Python (expert), C++, Java, C, SQL

**ML/AI Frameworks:** PyTorch, TensorFlow, Hugging Face Transformers, DeepSpeed

**Specialized Skills:** Mechanistic interpretability (SAEs, activation steering, neuron attribution), RLHF techniques (PPO, DPO, GRPO), Cross-lingual analysis

**Relevant Courses:** Advanced NLP, Advanced Topics in Security, Privacy, and Machine Learning, User-centered ML, LLMs Post Pretraining.

## TEACHING EXPERIENCE AND EXTRACURRICULAR

- **Teaching Assistant, CS 105: Introduction to Computing** (Fall 2023, Spring 2024, Fall 2024, Spring 2025). Assisted in course delivery, grading, and conducting lab sessions; awarded **Outstanding TA** for Spring 2024.
- Volunteer, **EMNLP 2021** and **NAACL 2022**. Supported conference logistics and session coordination.
- Managing Committee Member, **IEEE Student Branch**. Organized technical events and coordinated student outreach activities.