

KEYU HE

Los Angeles, CA — (213) 713-2973 — frankhe@usc.edu — [keyu-he.github.io](https://github.com/keyu-he)

Education

University of Southern California, Los Angeles, CA

August 2021 – May 2025

Double Major in Computer Science and Applied and Computational Mathematics

GPA: 3.98/4.00

Minor in Artificial Intelligence Applications

Dean's List Fall 2021, Spring 2022, Fall 2022, Spring 2023, Fall 2023, Spring 2024

Member of Phi Kappa Phi Honor Society

USC Center for Undergraduate Research in Viterbi Engineering (CURVE) Fellow

Research Interests

Explainable NLP Systems, Interpretable Machine Learning, Cross-Disciplinary AI Research, etc.

Relevant Coursework

Language Models in Natural Language Processing (A)

Applied Machine Learning for Natural Language Processing (A)

Directed Research (A)

Introduction to Machine Learning (A)

Applications of Machine Learning (A)

Professional C++ (A)

Introduction to Artificial Intelligence (A)

Applied Neural Networks (A)

Introduction to Data Analysis (A)

Introduction to Computer Systems (A)

Introduction to Embedded Systems (A)

Introduction to Internetworking (A)

Data Structure and Object-Oriented Design (A)

Introduction to Algorithms and the Theory of Computing (A-)

Linear Algebra and Differential Equations (A)

Probability Theory (A)

Mathematical Statistics (A)

Numerical Methods (A)

Conference Publications and Working Papers

- Huihan Li*, Arnav Goel*, **Keyu He**, and Xiang Ren. *Attributing Culture-Conditioned Generations to Pretraining Corpora*. ICLR 2025. [10.48550/arXiv.2412.20760](https://arxiv.org/abs/2412.20760).
- Brihi Joshi*, **Keyu He***, Kaitlyn Zhou, Sadra Sabouri Halestani, Souti Chattopadhyay, Swabha Swayamdipta, Xiang Ren. *Assessing Language Models' Capability to Explain to Different Audiences*. Under preparation. Aiming for ACL 2025.
- Keyu He**, Brihi Joshi, Tejas Srinivasan, Swabha Swayamdipta. *Beyond the Text: How Explanation Qualities Influence User Trust in Visual Language Models*. Under preparation. Aiming for ACL 2025.

(* Indicates equal contribution)

Research Experience

Center for Undergraduate Research in Viterbi Engineering (CURVE) Fellow

— USC Viterbi School of Engineering, US

Jan. 2024 - Present

- Contributing to the “Assessing Language Models' Capability to Explain to Different Audiences” project under Prof. Xiang Ren, Prof. Swabha Swayamdipta and PhD mentor Brihi Joshi, focusing on the comparative study of mechanistic and teleological explanations of “Why” questions by humans and GPT-4.
- Conducted quantitative analysis using lexical and semantic metrics and qualitative assessment to evaluate text complexity and alignment with user roles in curiosity-driven scenarios.
- Engaged in interdisciplinary workshops on cognitive science and NLP to refine methodologies for assessing context-driven explanation types.
- Fellowship awarded multiple times (Spring 2024, Summer 2024, Fall 2024, and Spring 2025), with a total funding amount of \$6,750.

Research Contributor on Visual Language Models (VLM) Project

— USC Viterbi School of Engineering, US

Jun. 2024 - Present

- Working on the “Beyond the Text: How Explanation Qualities Influence User Trust in Visual Language Models” project under Prof. Swabha Swayamdipta and PhD mentors Brihi Joshi and Tejas Srinivasan, focusing on developing tools to assess the faithfulness, relevance, and completeness of VLM rationales, aiming for helping users better judge whether to rely on explanations.
- Conducting both automatic and human evaluations to identify limitations of current text-only metrics and explore new vision-specific metrics.

Projects

LLM Prompt Recovery Project — USC

Mar. 2024 - Apr. 2024

- Developed a system to recover user prompts given original text and modified text generated by Gemma.
- Fine-tuned the Mixtral model using custom metrics, achieving a score of 0.65 with sentence-T5-base and sharpened cosine similarity (exponent = 3).
- Awarded a silver medal in the Kaggle competition for outstanding performance (ranked 75/2175, top 3.4%).
- See the final fine-tuned model here: [Mixtral-8x7b Instruct Finetuned](#).

AI-Based Career Advisor — USC

Nov. 2024 - Dec. 2024

- Developed an AI advisor to assist users in planning career paths based on skills and interests, leveraging datasets such as the JobSkills Dataset (1.3M entries) and LinkedIn Jobs Dataset.
- Designed a Streamlit-based interactive UI for user input and career recommendations, integrating GPT-4o for skill-based job suggestions and resource recommendations.
- Implemented a cosine similarity search on sentence embeddings for matching user skills with most-fit jobs and identifying skill gaps.
- Integrated Bing AI for real-time resource and job application link retrieval, enhancing usability with advanced support metrics for internal consistency verification.
- Enabled post-hoc evaluation using a T5 fine-tuned entailment verification model to validate skill-job relevance, ensuring reliable recommendations.

Enhancing Debugging Skills of LLMs with Prompt Engineering — USC

Aug. 2023 – Nov. 2023

- Improved the debugging capabilities of LLMs using innovative prompt engineering techniques.
- Conducted experiments using various prompting strategies (Zero-Shot, Few-Shot, Chain of Thought) to enhance the efficiency of GPT models in debugging tasks.

Automated Hate Speech Detection in Social Media — USC

Sep. 2023 – Dec. 2023

- Led the development of an advanced machine learning model for detecting hate speech on social media, employing a mix of techniques with a focus on BERT fine-tuning.
- Achieved a 94% accuracy rate in classification tasks, underlining the model's effectiveness in enhancing online safety and inclusivity through rigorous evaluation and optimization strategies.

USC Study Room/Area Rating Web Application — USC

Aug. 2022 – Dec. 2022

- Contributed to a team in the development of a web interface that allows students to rate and comment on study rooms located throughout the USC campus.
- Implemented sorting and filtering features, as well as a simplified reservation system, showcasing technical skills and contribution to the application's functionality.

Teaching Experience

Teaching and Grading Assistant — USC, US

Sep. 2022 - Present

- Selected for multiple roles: **Course Producer** for CSCI-102 and CSCI-360, Grader for MATH-117, MATH-126, MATH-129 and MATH-226.
- Ensured a consistent approach to teaching and grading by regularly collaborating with faculty and fellow graders.

Major Awards

- Silver Medal, Kaggle Competition, Ranked 75/2175 (Top 3.4%) on the global leaderboard, LLM-Prompt-Recovery Project, 2024
- USC Academic Achievement Award, Fall 2022, Spring 2023, Spring 2024, Fall 2024.
 - This award covered 11 units of tuition costs in total, amounting to approximately \$24,000.
- 4th Place, USC Integral Bee Competition, 2022
- 1st Prize, International Linguistics Olympiad (Senior Level), Individual Open Round, China, 2021
- 1st Prize, International Linguistics Olympiad (Senior Level), Team Open Round, China, 2021

Technical Skills

Programming: C++, C, Python, Java, MySQL, HTML, CSS, JS

Software: L^AT_EX, Mathematica, Matlab

Language: Mandarin (native), English (professional)