

Kellen Cheng

✉ kellantan@princeton.edu 4074094268 🏠 Website 📧 kellantan 🇺🇸 U.S. Citizen

Education

Princeton University

PH.D. ELECTRICAL AND COMPUTER ENGINEERING, GPA: 4.00

Aug. 2022 - Present

PRINCETON, NJ

Princeton University

M.A. ELECTRICAL AND COMPUTER ENGINEERING, GPA: 4.000

Aug. 2022 - Nov. 2024

PRINCETON, NJ

University of California, Los Angeles (UCLA)

B.S. ELECTRICAL ENGINEERING, TECHNICAL BREADTH COMPUTER SCIENCE, GPA: 3.926

Sept. 2018 - Jun. 2022

LOS ANGELES, CA

Experience

IBM Research: AI Research Scientist Intern

MENTOR: DR. ANNA LISA GENTILE

May. 2025 - Aug. 2025

SAN JOSE, CA

- Dynamic factuality assessment and evaluation through incorrect and misleading assertions (Current).

Samsung Research America: NLP Research Scientist Intern

MENTOR: DR. GANESH RAMESH

Jan. 2025 - May. 2025

MOUNTAIN VIEW, CA

- Designed a multi-agentic LLM framework to improve conversation summarization with iterative text feedback.
- Adapted the framework for on-device local inference with Apple Silicon using Ollama and MLX-LM.
- Performed quantized LoRA (QLoRA) instruction fine-tuning for a range of language models, from 0.5B to 9B parameters.

IBM Research: AI Research Scientist Intern

MENTOR: DR. ANNA LISA GENTILE

Jun. 2024 - Sept. 2024

SAN JOSE, CA

- Synthesized and curated an evaluation benchmark for health advice guardrails from Common Crawl web text.
- Designed and implemented a sparse human-in-the-loop system for semi-automatic annotation of synthetic data at scale.
- Formulated a method to automatically generate synthetic data using compact LLMs for health advice guardrails.
- Fine-tuned scalable and compact detector models on a blend of synthetic and open-source training data, beating GPT-4o by 3.73% in accuracy and 1.54% in F1-score, despite containing 400x less parameters.
- Created an internal Rest API that integrated my detector model and automated internal model evaluations for the team.
- Work published in the EMNLP industry track (first-author), with another work currently in submission (first-author).
- Filed a patent (stage 2) detailing a continual learning framework with model version-control and knowledge distillation for AI safety guardrail detector development.

Princeton: NLP Researcher

ADVISOR: DR. SUMA BHAT

Nov. 2022 - Nov. 2024

PRINCETON, NJ

- Created an end-to-end two-mask infilling fine-tuning objective for idiomatic knowledge injection using the IEKG dataset.
- Implemented two-stage fine-tuning with transfer learning to achieve new state-of-the-art performance of 83.75% accuracy on the IMPLI benchmark, an improvement of 12% compared to previous state-of-the-art.
- Conducted ablation and data perturbation studies to gauge contextual reasoning capabilities for off-the-shelf language models ranging from 0.5B to 7B parameters, uncovering that they actually perform *better* without the context.
- Work published in the NAACL (first-author) and EMNLP (second-author) main conferences.

Publications

- **Kellen Tan Cheng**, Anna Lisa Gentile, Pengyuan Li, Chad DeLuca, Guang-Jie Ren. *Don't Be My Doctor! Recognizing Healthcare Advice in Large Language Models*. EMNLP 2024 Industry Track.
- **Kellen Tan Cheng**, Suma Bhat. *No Context Needed: Contextual Quandary In Idiomatic Reasoning With Pre-Trained Language Models*. NAACL 2024 Main.
- Ziheng Zeng, **Kellen Tan Cheng**, Srihari Venkat Nanniyur, Jianing Zhou, Suma Bhat. *IEKG: A Commonsense Knowledge Graph for Idiomatic Expressions*. EMNLP 2023 Main.
- **Kellen Tan Cheng**, Kunakorn Atchaneeyasakul, Zeid Barakat, David Liebeskind, Fabien Scalzo. *CT Perfusion Imaging of the Brain with Machine Learning*. ISVC 2021.

Awards & Organizations

| | |
|---------------------------------------|----------------|
| Toby & Jack Wolf Travel Grant | 2024 |
| Bede Liu Travel Grant | 2023 |
| Princeton ECE Departmental Fellowship | 2022 |
| Tau Beta Pi | 2020 - Present |
| IEEE Eta Kappa Nu (HKN) | 2019 - Present |
| UCLA Dean's Honor List | 2019 - 2022 |

Skills

| | |
|------------|---|
| Languages | Python, C++, MATLAB |
| Frameworks | PyTorch, Tensorflow, MLX, Transformers |
| Tools | Ollama, Slurm, AWS EC2, Anaconda/Mamba, Jupyter, LaTeX, MS Office |