Ting-Wen (Kathy) Ko

RESEARCH INTEREST

Causal reasoning, trustworthy AI, AI safety, model interpretability

EDUCATION

 M.Sc., Computational Statistics and Machine Learning, University College London (in London, UK) Courses: Probabilistic and Unsupervised Learning, Supervised Learning, Bayesian Deep Learning, Statistical NOPen-Endedness and General Intelligence, Statistical NLP Supervised by Mengyue Yang and Jun Wang. 	Sep 2024 — Sep 2025 Model and Analysis,
 M.Sc., Computer Science, National Taiwan University (in Taipei, Taiwan) GPA: 4.17/4.3 Supervised by Pu-Jen Cheng and Jyun-Yu Jiang. Thesis: Enhancing Retrieval Augmented Generation with Passage Combination. 	Sep 2022 — Jun 2024
 B.Sc., Undergraduate Honors Program of Electrical Engineering and Computer Science National Yang Ming Chiao Tung University (Previously National Chiao Tung University) (in Hsinchu, Taiwan) Overall GPA: 3.93/4.3 	Sep 2018 - Aug 2022
Research and Work Experience	
University College London Master's Thesis Research	Apr 2025 — Current
 Exploring spurious correlation in LLMs from the perspective of OOD detection. 	
MediaTek Research UK Deep Learning Intern	Jun 2024 — Aug 2024
• Researched on the model representation when LLM reasons, layer optimization through adaptive pruning, an techniques to improve transformer-based language models.	nd diffusion model
National Taiwan University	Sep 2023 — Jun 2024
Master's Thesis Research (Link: 🔗) (Code: 🗟)	
Developed an end-to-end trainable dense passage retriever that optimizes passage combination selection for QA with LLM	r retrieval-augmented
• Implemented comprehensive baseline systems, including sparse/dense retrieval and reranking methods	
Global Media Team, Yahoo! Project Intern (Code: 🗟)	Sep 2023 — Dec 2023
 Researched on multitask learning approaches for text readability assessment. Implemented unsupervised learning-to-rank to predict a comprehensive readability signal. 	
Global Media Team, Yahoo! Research Engineering Intern	Jul 2023 — Aug 2023
 Designed an internal LLM chatbot with vector database integration for Yahoo news articles. Applied retrieval-augmented generation methods for improved explainability and diverse search result by masearch and asymmetric embeddings. 	ax marginal relevance
Selected Projects	
 Deconstructing LLM Faithfulness (Python) Revisited existing faithfulness tests from a causal framework. Investigate the application of RL techniques to train language models using faithfulness as a reward metric. Research areas: causality, RL fine-tuning, AI safety, trustworthy AI 	Feb 2025 — Current
BANDITPROMPT: Steering Creativity in Image Generation (Python / LLMs / Diffusion Models)	Feb 2025 — Current

BANDITPROMPT: Steering Creativity in Image Generation (Python / LLMs / Diffusion Models)

- Developed a novelty search-based approach for generating diverse sets of images using LLM prompt optimization. •
- Implemented a Beta-Thompson multi-armed bandit algorithm to adaptively select specialized mutation strategies. •
- Conducted comprehensive evaluations demonstrating superior diversity metrics compared to existing baselines. •
- Research areas: Generative AI, evolutionary algorithms, multi-armed bandits, prompt optimization •

TEACHING EXPERIENCE

System Programming Teaching Assistant (C)

National Taiwan University

- Designed and implemented a simulation of a context switch system utilizing non-local jumps and signals in a class assignment.
- Conducted TA sessions, providing guidance and support to 100+ students on academic coursework.

SERVICE

Student Volunteer

The 46th International ACM SIGIR Conference

Awards & Honors

- 2024
 Phi Tau Phi Scholastic Honor Society of the Republic of China Honorary Membership

 2022
 Ministry of Education (MOE) Overseas Exchange Student Financial Assistance Grant (US \$4.6K)

 SKILLS
 Large Language Models, Retrieval-Augmented Generation, Language Modeling, Transformer Architectures, Causal Inference, Evolutionary Algorithms
- chitectures, Causal Inference, Evolutionary AlgorithmsProgrammingPython, PyTorch, HuggingFace, C/C++, Java, SQLToolkitsGit, Linux, Pandas, NumPy, Scikit-learn, NLTK, Streamlit

Last updated April 23, 2025

Sep 2022 – Dec 2022

Jul 2023