

Jinqiang Yu

Chatswood NSW 2067, Australia

✉ jinqiang.yuuu@gmail.com | [in linkedin.com/in/jinqiang-yu-404bb0187](https://www.linkedin.com/in/jinqiang-yu-404bb0187) | 🎓 Jinqiang Yu

Summary

I am passionate about machine learning and AI areas, with experience in the fields of machine/deep learning, NLP, CV and explainable AI for ML, in particular LLMs. With publications as the lead author in renowned AI conferences and journals like AAAI and JAIR, I have also served as a reviewer for top AI conferences such as AAAI and IJCAI. Explainable AI can be applied in machine learning areas, such as debugging ML models, improving LLM models through techniques like PEFT and instruction fine-tuning. Additionally, it can be applied to different domains, such as healthcare AI/ML systems, where it can identify disease-contributing elements in X-ray images. Stay updated on cutting-edge AI techniques like generative AI and large language models, I am expecting to leverage my technical skills to contribute to AI/ML projects.

Education

Monash University

Melbourne, Australia

PhD in Artificial Intelligence

Feb 2021 - Apr 2024

- **Thesis Topic:** Explainable AI for Machine Learning
- **Thesis Description:** The PhD project aims to deal with explainable AI and machine learning problems, including training interpretable machine learning (ML) models, generating accurate and concise explanations to explain ML models in various domains like NLP, CV, healthcare and general classification tasks, and applying explainable AI for ML, in particular LLMs, such as improving LLM performance and enabling LLMs to generate explainable outputs.

Monash University

Melbourne, Australia

Master of Information Technology

Mar 2019 - Dec 2020

- Graduated with H1
- **Minor Thesis Topic:** Computing optimal interpretable machine learning models.
- **Thesis Description:** The thesis focuses on interpretable models, aiming at developing advanced approaches to computing machine learning models that are both accurate and interpretable.

Experience

OPTIMA

Melbourne, Australia

PhD Researcher

Apr 2021 - Present

- I am engaged in a research project at OPTIMA, which focuses on explainable AI for machine learning, including generating interpretable ML models and computing accurate and concise explanations, aiming at developing methods to help users understand and explain ML model inferences and predictions, and applying explainable AI to improve ML performance in diverse domains, especially LLMs.

Monash University

Melbourne, Australia

Teaching Associate

Jun 2021 - Nov 2021

- Tutoring and grading students in tutorials, assignments, and final exams.
- **Unit:** FIT5220 - Solving discrete optimisation problems.
- **Topics:** Constraint Programming, Mixed Integer Programming, Boolean Satisfiability (SAT) Solving, Large Neighbourhood Search.

Professional Research Projects

Applying Explainable AI in Machine Learning and LLMs

Oct 2023 - present

- Conducting research on LLMs with the use of explainable AI. Aiming to improve LLMs' performance, help others understand the models' behavior, and build trust of the models.
- Developing approaches to improving the performance of LLMs with the use of Parameter-Efficient Fine-Tuning (PEFT) and Prompt Engineering given explanations.
- Developing methods to generate more explainable outputs of LLMs for different users with the use of Instruction Fine-Tuning and Reinforcement Learning with Human Feedback (RLHF).
- **Expertise:** Machine Learning, LLMs, Transformers, Generative AI, Prompting Engineering, PEFT, Instruction Fine-Tuning, RLHF, Explainable AI.

Explainability and Its Application in NLP and CV

Dec 2022 - present

- Conducting research on explainability in NLP and CV domains.
- Developing approaches to generating feature-attribution and feature selection based explanations in NLP/CV domains.
- Applying explainable AI methods for healthcare images.
- **Expertise:** LLMs, Transformers, CNNs, Natural Language Processing, Computer Vision, Explainable AI.

Computing Succinct and Accurate Explanations for ML

Feb 2021 - Apr 2023

- Conducted research on generating more accurate and concise post-hoc (*abductive*) explanations answering why the prediction is made and (*contrastive*) explanations targeting why not the prediction is made.
- Developed approaches to mining background knowledge and applying them to generating accurate and concise explanations in diverse domains.
- Developed approaches to generating feature attribution indicating the contribution of a feature.
- **Expertise: Machine Learning, Data Mining, Explainable AI.**

Learning Optimal Interpretable Machine Learning Models

Feb 2020 - Mar 2023

- Conducted research on developing approaches to training and computing interpretable machine learning models that can be easily understood by users.
- **Expertise: Machine Learning, Interpretable Models, Reasoning, Explainable AI.**

Publications

Anytime Approximate Formal Feature Attribution

Jinqiang Yu Graham Farr, Alexey Ignatiev, Peter J. Stuckey

arXiv preprint arXiv:2312.06973, 2023

On Formal Feature Attribution and Its Approximation

Jinqiang Yu, Alexey Ignatiev, Peter J. Stuckey

arXiv preprint arXiv:2307.03380, 2023

From Formal Boosted Tree Explanations to Interpretable Rule Sets

Jinqiang Yu, Alexey Ignatiev, Peter J. Stuckey

29th International Conference on Principles and Practice of Constraint Programming (CP), 2023

Eliminating the Impossible, Whatever Remains Must Be True: On Extracting and Applying Background Knowledge in the Context of Formal Explanations

Jinqiang Yu, Alexey Ignatiev, Peter J. Stuckey, Nina Narodytska, Joao Marques-Silva

37th AAAI Conference on Artificial Intelligence (AAAI), 2023

Learning Optimal Decision Sets and Lists with SAT

Jinqiang Yu, Alexey Ignatiev, Peter J. Stuckey, Pierre Le Bodic

Journal of Artificial Intelligence Research (JAIR) 72 (2021) pp. 1251–1279. 2021

Computing Optimal Decision Sets with SAT

Jinqiang Yu, Alexey Ignatiev, Peter J. Stuckey, Pierre Le Bodic

26th International Conference on Principles and Practice of Constraint Programming (CP), 2020

Skills

Programming Languages

Python, Java, R, C++, MATLAB

DS/ML/AI Technologies

Domains - NLP, CV, LLM, Generative AI, Explainable AI

ML/DL/GenAI Models - LLM, Transformer, CNN, Diffusion model, GAN, VAE, etc

Modelling - PyTorch, Tensorflow, Scikit-learn, SciPy, Hugging Face

LLM Technologies - LLM framework (LangChain), RAG, Parameter-Efficient Fine-Tuning (PEFT), Instruction Fine-Tuning, RLHF, and Prompt Engineering

Cleansing and Wrangling - Pandas, Numpy, NLTK

Visualization - Matplotlib, Seaborn

Big Data Management - MongoDB, SQL, Spark

Miscellaneous Technologies

Github, Docker, Slurm, Linux

Awards and Scholarships

2021-2024 **Monash Graduate Scholarship**

Scholarship covers living expenses and tuition

2020 **Best Paper Award**

Our paper “Computing Optimal Decision Sets with SAT” has been selected for the Best Paper Award for the CP/ML Track of CP 2020.

Scientific Activities

2024 PC member of the AAAI Conference on Artificial Intelligence(AAAI-2024)

2024 PC member of the International Joint Conferences on Artificial Intelligence (IJCAI-2024)

References available upon request.