

Qian Chen

qianchen901005@gmail.com | github.com/qc31415926 | Tel: (+86)19370575316

EDUCATION

East China Normal University 985

Master of Computer Science and Technology

Hefei University of Technology 211

Bachelor of Computer Science and Technology

Sep. 2022 – June. 2025

Aug. 2018 – June 2022

RESEARCH EXPERIENCE

My research interests include Natural Language Processing, Interpretability and Large Language Models.

CIDR: A Cooperative Integrated Dynamic Refining Method for Minimal Feature Removal Problem

AAAI2024 First Author

- Minimal feature removal problem aims to find the minimum feature set. Previous works rely on monotonic assumptions, which cannot be satisfied in general scenarios.
- Prove that using Integrated Gradients we can transform the original problem into a knapsack problem and propose a plug-and-play method for generating minimum feature candidate sets.
- Extensive evaluation on Eraser benchmarks demonstrate the effectiveness of our method.

PE: A Poincare Explanation Method for Text Hierarchy Generation

EMNLP2024 Finding oa:3.25 meta:4

- To model non-contagious feature interactions, prior studies build the hierarchical attribution tree by enumerating all combinations, neglecting underlying syntax and semantics.
- Propose a hyperbolic probing method and introduce a fast algorithm for generating hierarchical attribution trees.
- Experimental results show the effectiveness of our approach in building high-quality hierarchical explanations.

Lottery Ticket Mask for Safety in LLMs Model Merge

AAAI2025 under review

- Prior methods ignore that merging different domain sft LLMs would cause safety performance downgrade.
- Propose a plug-and-play method to generate a mask and dynamically to pick the parameter to merge.
- Experimental results show our method can integrate existing model merge methods without losing performance on relevant domains.

Concept Based Continuous Prompts for Interpretable Text Classification

AAAI2025 under review

- Word-based explanation methods lacks comprehensive semantic understanding and cannot generalize to BPE encoding based models.
- Propose a framework to decompose continuous prompts to human-readable concepts.
- Experimental results show our framework can achieve the same results as the original P-tuning and word-based methods with a few concepts while providing more plausible results.

INTERNSHIP

Shanghai AI Lab OpenTrust Lab | Research Intern

July 2024 – Now

- Mainly working on **LLM model merge**

AWARDS

National Scholarship	2024.11
The First Prize Scholarship (top 10%)	2021.11
The Second Prize Scholarship	2020.11

SKILLS

Languages: Mandarin (Native), English (IELTS:7)

Relevant Course: Natural Language Understanding(93/100), Machine Learning Basic(93/100), Linear Algebra(99/100), Probability Theory and Mathematical Statistics(98/100), Advanced Mathematics A(93/100), Advanced Mathematics B(95/100)