

Zhuofeng Wu

Applied Scientist II at Amazon foundational LLM - Rufus

Ph.D. from School of Information, University of Michigan

(E) zhuofeng@umich.edu (C) (734) 780-9664 (W) <https://cserxy.github.io/>

EDUCATION

University of Michigan, Ann Arbor, US

Aug 2018 – Apr 2024

Ph.D. from School of Information

Advisor: V.G. Vinod Vydiswaran

Committee: Rada Mihalcea, Chaowei Xiao, Paramveer Dhillon

Zhejiang University, Hangzhou, China

Sept 2013 - Jun 2017

B.E. in the College of Computer Science & Chu Kochen Honors College

Advisor: Fei Wu

Received waiver for the National College Entrance Exam to enter Zhejiang University from **1st Prize in National Olympiad in Informatics in Provinces**

WORKING EXPERIENCE

Amazon, foundational LLM team - Rufus

May 2024 – present

Applied Scientist II

Manager: Huasheng Li

Apple, Machine Learning Research team

Apr 2023 – Aug 2023

ML Research Intern

Mentors: Yizhe Zhang and Navdeep Jaitly

Meta, AI Integrity team

May 2021 – Aug 2021

Research Intern

Mentors: Sinong Wang and Hao Ma

Meta, AI Integrity team

May 2020 – Aug 2020

Research Intern

Mentors: Sinong Wang and Hao Ma

Alibaba Group

May 2019 – Aug 2019

Research Intern

Mentor: Fei Sun

PUBLICATIONS

Divide-or-Conquer? Which Part Should You Distill Your LLM?

Zhuofeng Wu, He Bai, Aonan Zhang, Jiatao Gu, VG Vydiswaran, Navdeep Jaitly, Yizhe Zhang

In Proceedings of Findings EMNLP 2024. ([pdf](#))

HiCL: Hierarchical Contrastive Learning of Unsupervised Sentence Embeddings

Zhuofeng Wu, Chaowei Xiao, V. G. Vydiswaran

In Proceedings of Findings EMNLP 2023. ([pdf](#))

Improving Large Language Models Function Calling and Interpretability via Guided-Structured Templates

Hy Dang, Tianyi Liu, Zhuofeng Wu, Jingfeng Yang, Haoming Jiang, Tao Yang, Pei Chen, Zhengyang Wang, Helen

Wang, Huasheng Li, Bing Yin, Meng Jiang

In Proceedings of EMNLP 2025. ([pdf](#))

UniConv: Unifying Retrieval and Response Generation for Large Language Model in Conversation

Fengran Mo, Yifan Gao, Chuan Meng, Xin Liu, Zhuofeng Wu, Kelong Mao, Zhengyang Wang, Pei Chen, Zheng Li,

Xian Li, Bing Yin, Meng Jiang

In Proceedings of ACL 2025. ([pdf](#))

IDPG: An Instance-Dependent Prompt Generation Method

Zhuofeng Wu, Sinong Wang, Jiatao Gu, Rui Hou, Yuxiao Dong, V. G. Vydiswaran, Hao Ma

In Proceedings of NAACL 2022 (Oral Presentation). ([pdf](#)) ([video](#))

PLANNER: Generating Diversified Paragraph via Latent Language Diffusion Model

Yizhe Zhang, Jiatao Gu, Zhuofeng Wu, Shuangfei Zhai, Josh Susskind, Navdeep Jaitly
In Proceedings of NeurIPS 2023. ([pdf](#))

Leveraging historical information to boost retrieval-augmented generation in conversations

Fengran Mo, Yifan Gao, Zhuofeng Wu, Xin Liu, Pei Chen, Zheng Li, Zhengyang Wang, Xian Li, Bing Yin, Meng Jiang, Jian-Yun Nie

In Information Processing & Management Journal. ([pdf](#))

Defending against Insertion-based Textual Backdoor Attacks via Attribution

Jiazhao Li, Zhuofeng Wu, Wei Ping, Chaowei Xiao, V. G. Vydiswaran
In Proceedings of Findings ACL 2023. ([pdf](#))

Chatgpt as an attack tool: Stealthy textual backdoor attack via blackbox generative model trigger

Jiazhao Li, Yijin Yang, Zhuofeng Wu, V. G. Vydiswaran, Chaowei Xiao
In Proceedings of NAACL 2024. ([pdf](#))

Identify Shifts of Word Semantics through Bayesian Surprise

Zhuofeng Wu, Cheng Li, Zhe Zhao, Fei Wu, Qiaozhu Mei
In Proceedings of SIGIR 2018 (Oral Presentation). ([pdf](#))

PREPRINT

Clear: Contrastive learning for sentence representation

Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa, Fei Sun, Hao Ma
arXiv preprint arXiv:2012.15466 (2020). ([pdf](#))

Adversarial Demonstration Attacks on Large Language Models

Jiongxiao Wang, Zichen Liu, Keun Hee Park, Zhuojun Jiang, Zhaocheng Zheng, Zhuofeng Wu, Muhan Chen, Chaowei Xiao

arXiv preprint arXiv:2305.14950 ([pdf](#))

PROJECT EXPERIENCE

Amazon, Shopping LLM - Rufus

May 2024 – present

[RAG]: Retrieval Augmented Generation Improvement via Context Denoising

- Enhanced LLM accuracy by implementing conversation history denoising for Amazon product detail pages.
- Designed and optimized pipeline achieving 70% win-tie rate against baseline performance and reducing hallucination.

[Agent]: Product Recommendation System Integration in LLM

- Pioneered and launched first-ever product recommendation system for Rufus platform via agent.
- Led end-to-end development of post-training product selection optimization using Direct Preference Optimization (DPO).

[Scaling Law]: Estimated the Target Model's Hyperparameters

- Built and trained a set of sampling models to estimate the optimized batch size, learning rate for the target model.

[Reasoning & Personalization]: Integrate Personalization into Reasoning/Thinking LLM.

- Developed multi-source context integration framework to enhance LLM's personalization capabilities using Proximal Policy Optimization (PPO) and Group Relative Policy Optimization (GRPO).

Apple, Machine Learning Research

Apr 2023 – Aug 2023

[Reasoning]: Knowledge Distillation from LLM to Small Models: A Perspective from Question Decomposition

- Pioneered a novel knowledge distillation approach that leverages GPT-4 to decompose complex questions into structured sub-questions, enabling smaller models to adopt advanced reasoning capabilities.
- Developed a two-stage training pipeline that combines initial fine-tuning on question-subquestion pairs with reinforcement learning (PPO), using GPT-4-generated rewards to optimize performance.
- Showed that distilling question decomposition (not solving) into small models yields GPT-level planning at a fraction of the cost, improving GSM8K/DROP performance and sometimes surpassing the teacher.
- This work was presented at EMNLP'24 as poster.

Meta, AI Integrity

May 2021 – Aug 2021

[Parameter-Efficient Fine-Tuning]: IDPG: An Instance-Dependent Prompt Generation Method

- First customized prompt for each input rather than one prompt for all inputs.
- Offered comparable performance to Adapter-based methods while using fewer parameters.
- Extensive evaluations on ten natural language understanding tasks show that IDPG consistently outperforms task-specific prompt tuning methods by 1.6–3.1 points.

- This work was presented at **NAACL'22** as **oral**.

Meta, AI Integrity

May 2020 – Aug 2020

[*Pre-training*]: CLEAR: Contrastive Learning for Sentence Representation

- Proposed to align the representation of different argumentation for same sentence.
- Explored several argumentations and their combinations in the text domain.
- Revealed that different argumentations in pre-training enhance the model's different abilities.
- Outperformed several baselines (including BERT & RoBERTa) on GLUE & SentEval benchmark.

Alibaba Group

May 2019 – Aug 2019

[*Pre-training*]: Seg-BERT: A Hierarchical Structure for Document Classification

- Applied a hierarchical structure for the long text classification.
- Outperformed the state-of-the-art by a large margin on IMDB.
- Proposed to mask sentence in pre-training to improve the performance.

School of Information, University of Michigan

Apr 2016 – Apr 2018

Advisor: Prof. Qiaozhu Mei

[*Word Embedding*]: Identify Shifts of Word Semantics through Bayesian Surprise

- Explicitly established the stable topological structure of word semantics and identified the surprising changes over time.
- Proposed a statistical framework to apply **Bayesian Surprise** in detecting the meaning-changed words in **temporal-based word semantic networks**. This framework can be generalized to finding the change points in many other networks.
- Conducted experiments on ACM DL, DBLP and Google Books Ngram data set for synthetic evaluation which artificially introducing changes to a corpus. Outperformed the state-of-the-art by a large margin.
- This work was presented at **SIGIR'18** as **oral** and was adopted as a part of a **KDD'18 Workshop Keynote Talk** “Identifying Shifts in Evolutionary Semantic Spaces”.

[*Dimensionality Reduction*]: A Tool to Visualize the Evolution of Conference Topics

- Visualized a 40-year evolution of data science related communities and embedded papers, keywords, authors in the same space.
- Provided a powerful tool for researchers to model the research focus of different conferences.
- This work was presented in an invited talk in **KDD'18 Deep Learning Day** by Prof. Mei.

SERVICE

Area Chair

ACL 2025, EMNLP 2025

Student Volunteer

NAACL 2022, SIGIR 2018

Conference Reviewer

ICLR 2026, NeurIPS 2025, ICML 2025, AAAI 2025, EMNLP 2024, ICML 2024, ICLR 2024, NeurIPS 2023, EMNLP 2023,
ACL 2023, SIGIR 2023

ACL Rolling Reviewer

October 2025 Cycle, June 2024 Cycle, April 2024 Cycle, February 2024 Cycle, December 2023 Cycle, October 2023 Cycle,
June 2023 Cycle, April 2023 Cycle, December 2022 Cycle

AWARDS

EMNLP Student Travel Grant from Big Picture Workshop, 2023.

SIGIR Student Travel Grant, 2018.

Outstanding Graduates of Zhejiang Province, 2017.

3rd Prize in Collegiate Programming Contest of Zhejiang University, 2014, 2015.

2nd Prize of Excellent Undergraduate Scholarship, 2014.

1st Prize in National Olympiad in Informatics in Provinces in 2012.

1st Prize in National Olympiad in Mathematics in Provinces in 2010.

MENTORING

Hy Dang, Applied Scientist Intern, PhD student from University of Notre Dame

Sept 2024 – May 2025

Hongye Jin, Applied Scientist Intern, now Applied Scientist II at Amazon Rufus

Sept 2024 – Dec 2024

Fengran Mo, Applied Scientist Intern, PhD student from Université de Montréal

May 2024 – Dec 2024

Tian Xia, Undergraduate student from University of Michigan

May 2023 – May 2024

Jiazhao Li, PhD student from University of Michigan, now Applied Scientist at Amazon Rufus

Sept 2021 – May 2024