

CONTACT INFORMATION	Department of Computer Science University of Virginia Room 224, Rice Hall Charlottesville, VA 22903 Email : zhepei.wei@virginia.edu Homepage: weizhepei.com	
RESEARCH INTERESTS	My recent research mainly focuses on large language model (LLM), including but not limited to retrieval-augmented generation (RAG), efficient inference, and alignment with human preference.	
EDUCATION	University of Virginia <ul style="list-style-type: none"> • Ph.D., Department of Computer Science • Advisor: Yu Meng 	Charlottesville, U.S. Aug 2022 - present
	Jilin University <ul style="list-style-type: none"> • M.S., School of Artificial Intelligence • B.S., College of Computer Science and Technology • Advisor: Yi Chang 	Changchun, China Sep 2019 - Jun 2022 Sep 2015 - Jul 2019
	University of Virginia <i>Research Assistant</i> <ul style="list-style-type: none"> • Advisor: Yu Meng • Large Language Model (LLM): Research on LLM alignment, retrieval-augmented generation, and efficient LLM inference. 	Charlottesville, U.S. Aug 2023 - present
	University of Virginia <i>Research Assistant · Human-Centric Data Mining Lab</i> <ul style="list-style-type: none"> • Advisor: Hongning Wang • Multi-agent Decision Making: Research on incentivizing reinforcement learning (RL) agents for collaborative decision-making in federated environment. 	Charlottesville, U.S. Aug 2022 - Aug 2023
	Jilin University <i>Research Assistant</i> <ul style="list-style-type: none"> • Advisor: Yi Chang • Language Understanding and Reasoning: Research on information extraction and natural language reasoning tasks with knowledge graph and pre-trained LMs. 	Changchun, China Sep 2019 - Jun 2022
	Jilin University <i>Research Assistant · Machine Learning Lab</i> <ul style="list-style-type: none"> • Advisor: You Zhou • Autonomous Agent: Research on improving deep Q-learning (reinforcement learning) for autonomous decision-making in an interactive game environment. 	Changchun, China May 2017 - Jun 2018
	Amazon <i>Applied Scientist Intern · Rufus Alignment Team</i> <ul style="list-style-type: none"> • Manager: Wenlin Yao, Lihong Li • Multi-modal Reasoning: Work on improving the reasoning ability of multi-modal LLMs by scaling up inference-time compute. 	Seattle, U.S. Jan 2025 - present
WORK EXPERIENCE		

SELECTED
PUBLICATIONS

Google Scholar Citations: 700+, Github Stars: 900+ (as of Jan. 2025)

1. **Zhepei Wei**, Wei-Lin Chen, Xinyu Zhu, Yu Meng. *Fast and Accurate Language Model Decoding via Parallel Token Processing*. Under Review. Presented at NeurIPS 2024 AFM Workshop (Oral: 8/157).
2. Shangjian Yin, **Zhepei Wei**, Xinyu Zhu, Wei-Lin Chen, Yu Meng. *Self-Alignment Optimization for Language Models*. Under Review.
3. **Zhepei Wei**, Wei-Lin Chen, Yu Meng. *InstructRAG: Instructing Retrieval-Augmented Generation via Self-Synthesized Rationales*. **ICLR 2025**.
4. **Zhepei Wei**, Chuanhao Li, Tianze Ren, Haifeng Xu, Hongning Wang. *Incentivized Truthful Communication for Federated Bandits*. **ICLR 2024**.
5. **Zhepei Wei**, Chuanhao Li, Haifeng Xu, Hongning Wang. *Incentivized Communication for Federated Bandits*. **NeurIPS 2023**.
6. Erxin Yu, Lan Du, Yuan Jin, **Zhepei Wei**, Yi Chang. *Learning Semantic Textual Similarity via Topic-informed Discrete Latent Variables*. **EMNLP 2022**.
7. **Zhepei Wei**, Yue Wang, Jinnan Li, Zhining Liu, Erxin Yu, Yuan Tian, Xin Wang, Yi Chang. *Towards Unified Representations of Knowledge Graph and Expert Rules for Machine Learning and Reasoning*. **AACL-IJCNLP 2022**.
8. Runliang Niu, **Zhepei Wei**, Yan Wang, Qi Wang. *AttExplainer: Explain Transformer via Attention by Reinforcement Learning*. **IJCAI-ECAI 2022**.
9. **Zhepei Wei**, Yue Wang, Jinnan Li, Yingqi Guo, Zhining Liu, Erxin Yu, Yi Chang. *CogKG: Unifying Symbolic and Continuous Knowledge Representations for Machine Reasoning*. AKBC 2021 USKB workshop.
10. **Zhepei Wei**, Jianlin Su, Yue Wang, Yuan Tian, Yi Chang. *A Novel Cascade Binary Tagging Framework for Relational Triple Extraction*. **ACL 2020**.

AWARDS &
HONORS

- Copenhagen Charitable Trust Bicentennial Fellowship (2024) UVA
- OpenAI Researcher Access Program Recipient (2024) OpenAI
- John A. Stankovic Graduate Research Award (2024) UVA
- NeurIPS Scholar Award (2023) NeurIPS Foundation
- Computer Science Scholar Fellowship (2022) University of Virginia (UVA)
- National Scholarship (2018, 2020) Ministry of Education of China
- Outstanding Graduate Award (2020) Jilin University
- Honor Graduation Award (2019) Jilin University
- Outstanding Undergraduate Award (2016, 2017, 2018) Jilin University
- National Endeavor Scholarship (2016) Ministry of Education of China

ACADEMIC
SERVICES

Conference Paper Reviewer

- The International Conference on Learning Representations (ICLR), 2025
- The International Conference on Machine Learning (ICML), 2024, 2025
- Annual Conference on Neural Information Processing Systems (NeurIPS), 2024
- The AAAI Conference on Artificial Intelligence (AAAI), 2020, 2021, 2024, 2025
- The Web Conference (WWW), 2023
- Conference on Empirical Methods in Natural Language Processing (EMNLP), 2021
- Annual Meeting of the Association for Computational Linguistics (ACL), 2020