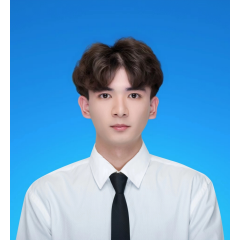


Ru Peng 彭儒

Homepage: pengr.github.io ◇ [Google Scholar](#) ◇ [Github](#)

Email: rupeng@zju.edu.cn ◇ Phone: (+86)13802728634



Research Interest and Highlights

My research interests spread across several AI fields, including **LLMs (current emphasis)**, Machine Learning, Natural Language Processing, Multimodal.

Large Language Models: focus on pre-training data management, including data selection, mixing, synthesis for LLMs.

Automated Model Evaluation: focus on unsupervised model evaluation on varied environments.

Machine Translation: focus on multi-modal (vision-language), sign, text-only machine translation.

Previously, I have conducted extensive research on the **full LLM pipeline**, covering both data and model sides. I maintain two popular GitHub repository in **LLM-Synthetic-Data (285+ stars)** and **TableGPT (590+ stars, 24.1k downloads last month)**. Google Scholar citations **2700+**.

Education

Zhejiang University

Ph.D. of Computer Science

Sep 2022 - Present

Advisor: Junbo Zhao and Gang Chen

Guangdong University of Technology

M.S. of Communication Engineering

Sep 2017 - Jul 2020

Advisors: Yi Fang and Tianyong Hao

Liaoning Technical University

B.E. of Communication Engineering

Sep 2013 - Jul 2017

Work Experience

Qwen Pre-training Team, Alibaba Group

Research Intern on Large Language Models

Oct 2023 - April 2025

Mentor: Dayiheng Liu, Junyang Lin and Chang Zhou

- Contributed to the **Qwen 1.5/2/2.5/3** series base models.
- Developed the **Data Manager (DataMan)** to select and mix data for pre-training large language model.
- Contributed to **data synthesis in open-ended tasks** for the Qwen series models.

Institute of Computer Innovation, Zhejiang University

Research Intern on Machine Translation and Automated Model Evaluation

Sep 2021 - Sep 2022

Mentor: Junbo Zhao and Peng Lu

ATTL Lab, NICT

Research Intern (Remote) on Machine Translation

Jul 2020 - Sep 2021

Mentor: Kehai Chen

Publications

[1] **Qwen3 technical report**. Qwen Team. Arxiv 2025.

[2] **Qwen2.5 technical report**. An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Biao Sun, Bin Luo, Bin Zhang, Binghai Wang, Chaojie Yang, Chang Si, Cheng Chen, Chengpeng Li, Chujie Zheng, Fan Hong, Guanting Dong, Guobin Zhao, Hangrui Hu, Hanyu Zhao, Hao Lin, Hao Xiang, Haoyan Huang, Humen Zhong, Jialin Wang, Jialong Tang, Jiandong Jiang, Jianqiang Wan, Jianxin Ma, Jianyuan Zeng, Jie Zhang, Jin Xu, Jinkai Wang, Jinzheng He, Jun Tang, Ke Yi, Keqin Chen, Langshi Chen, Le Jiang, Lei Zhang, Liang Chen, Man Yuan, Mingkun Yang, Minmin Sun, Na Ni, Nuo Chen, Peng Wang, Peng Zhu, Pengcheng Zhang, Pengfei Wang, Qiaoyu Tang, Qing Fu, Rong Zhang, **Ru Peng**,

Ruize Gao, Shanghaoran Quan, Shen Huang, Shuai Bai, Shuang Luo, Sibao Song, Song Chen, Tao He, Ting He, Wei Ding, Wei Liao, Weijia Xu, Wenbin Ge, Wenbiao Yin, Wenyuan Yu, Xianyan Jia, Xianzhong Shi, Xiaodong Deng, Xiaoming Huang, Ximing Zhou, Xinyu Wang, Xipin Wei, Xuejing Liu, Yang Liu, Yang Yao, Yang Zhang, Yibo Miao, Yidan Zhang, Yikai Zhu, Yinger Zhang, Yong Jiang, Yong Li, Yongan Yue, Yuanzhi Zhu, Yunfei Chu, Zekun Wang, Zhaohai Li, Zheren Fu, Zhi Li, Zhibo Yang, Zhifang Guo, Zhipeng Zhang, Zhiying Xu, Zile Qiao, Ziyi Meng. Qwen Team. Arxiv 2025.

[3] **Qwen2 technical report**. An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, **Ru Peng**, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhihao Fan. Arxiv 2024.

[4] **Introducing qwen1.5**. Qwen Team. February, 2024.

[5] DataMan: Data Manager for Pre-training Large Language Models.
Ru Peng, Kexin Yang, Yawen Zeng, Junyang Lin, Dayiheng Liu, Junbo Zhao.
The Thirteenth International Conference on Learning Representations (**ICLR**), 2025.

[6] LLM-Enhanced Query Generation and Retrieval Preservation for Task-Oriented Dialogue.
Jiale Chen, Xuelian Dong, Wenxiu Xie, **Ru Peng**, Kun Zeng, Tianyong Hao.
Findings of the Association for Computational Linguistics: **ACL** 2025.

[7] Predicting Rewards Alongside Tokens: Non-disruptive Parameter Insertion for Efficient Inference Intervention in Large Language Model.
Chenhan Yuan, Fei Huang, **Ru Peng**, Keming Lu, Bowen Yu, Chang Zhou, Jingren Zhou.
Conference on Empirical Methods in Natural Language Processing (**EMNLP**), 2024.

[8] Embedding and Gradient Say Wrong: A White-Box Method for Hallucination Detection.
Xiaomeng Hu, Yiming Zhang, **Ru Peng**, Haozhe Zhang, Chenwei Wu, Gang Chen, Junbo Zhao.
Conference on Empirical Methods in Natural Language Processing (**EMNLP**), 2024.

[9] Inference-Time Decontamination: Reusing Leaked Benchmarks for Large Language Model Evaluation.
Qin Zhu, Qingyuan Cheng, Runyu Peng, Xiaonan Li, Tengxiao Liu, **Ru Peng**, Xipeng Qiu, Xuanjing Huang.
Findings of the Association for Computational Linguistics: **EMNLP** 2024.

[10] DORY: Deliberative Prompt Recovery for LLM.
Lirong Gao, **Ru Peng**, Yiming Zhang, Junbo Zhao.
Findings of the Association for Computational Linguistics: **ACL** 2024.

[11] Energy-based Automated Model Evaluation.
Ru Peng, Heming Zou, Haobo Wang, Yawen Zeng, Zenan Huang, Junbo Zhao.
The Twelfth International Conference on Learning Representations (**ICLR**), 2024.

[12] CAME: Contrastive Automated Model Evaluation.
Ru Peng, Qiuyang Duan, Haobo Wang, Jiachen Ma, Yanbo Jiang, Yongjun Tu, Xiu Jiang, Junbo Zhao.
IEEE/CVF International Conference on Computer Vision (**ICCV**), 2023.

[13] Distill The Image to Nowhere: Inversion Knowledge Distillation for Multimodal Machine Translation.
Ru Peng, Yawen Zeng, Junbo Zhao.
Conference on Empirical Methods in Natural Language Processing (**EMNLP**), 2022. **(Oral, top 4.1%)**

[14] HybridVocab: Towards Multi-Modal Machine Translation via Multi-Aspect Alignment.
Ru Peng, Yawen Zeng, Junbo Zhao.
International Conference on Multimedia Retrieval (**ICMR**), 2022. **(Oral)**

[15] Deps-SAN: Neural Machine Translation with Dependency-Scaled Self-Attention Network.
Ru Peng, Nankai Lin, Yi Fang, Shengyi Jiang, Tianyong Hao, Boyu Chen, Junbo Zhao.
The 29th International Conference on Neural Information Processing (**ICONIP**), 2022. **(Oral)**

[16] Syntax-Aware Attentional Neural Machine Translation Directed by Syntactic Dependency Degree.
Ru Peng, Tianyong Hao, Yi Fang.
Neural Computing and Applications (**NCA**), 2020. **(Impact Factor: 5.6)**

- [17] Neural Machine Translation with Attention Based on a New Syntactic Branch Distance.
Ru Peng, Zhitao Chen, Tianyong Hao, Yi Fang.
The 15th China Conference on Machine Translation (CCMT), 2019. **(Best Paper Candidates 3th)**

Preprints

- [1] DotaMath: Decomposition of Thought with Code Assistance and Self-correction for Mathematical Reasoning.
Chengpeng Li, Guanting Dong, Mingfeng Xue, **Ru Peng**, Xiang Wang, Dayiheng Liu.
Arxiv 2024.
- [2] Better Sign Language Translation with Monolingual Data.
Ru Peng, Yawen Zeng, Junbo Zhao.
arXiv preprint arXiv:2304.10844, 2023.
- [3] Image classification prediction method, device, equipment and storage medium.
Qiuyang Duan, **Ru Peng**, Junbo Zhao, Yongjun Tu, Xiu Jiang.
C.N. Patent , 2023.

Awards

Merit Student, Zhejiang University, 2022 – 2024
Academic Scholarship, GDUT, 2017 – 2019
Outstanding Student Scholarship, LNTU, 2016

Professional Service

Conference Reviewer: NeurIPS 22/23/24, ICML 23/25, ICLR 24/25, CVPR 25, ICCV 23/25, ECCV 24, ACL 24, AISTATS 25, COLM 2024

Journal Reviewer: IEEE Transactions on Big Data (TBD), Transactions of Machine Learning Research (TMLR)

Publication Chair: International Conference on Natural Language Processing (ICNLP) 2025

Invited Talks

Feb 2025, *DataMan: Data Manager for Pre-training Large Language Models*. Synced (机器之心).
Nov 2024, *Energy-based Automated Model Evaluation*. Renmin University of China
March 2024, *Energy-based Automated Model Evaluation*. DataWhale

Skills

Programming Languages: Java, Python, C, C#, Matlab, Shell, Perl
Deep Learning Frameworks: Pytorch, Tensorflow, Theano, Keras, Ray
Machine Learning Libraries: LLama-Factory, Megatron, Verl, Vllm, NumPy, SciPy, Sklearn, Matplotlib, Pandas, NLTK, etc
AI Experience: Rich experience in theories, toolkits, codebases, experiments for AI fields