

# Jianzhu Yao

<https://yao-jz.github.io>

Email : [cnyaojz@gmail.com](mailto:cnyaojz@gmail.com)

Mobile : +86-139-3693-8133

## Education

### Tsinghua University

Sep 2019 – Present

*Bachelor of Engineering in Computer Science and Technology; GPA: 3.86/4.00*

*Beijing, China*

**Courses:** Linear Algebra(4.0), Object-Oriented Programming(4.0), Calculus(4.0), Discrete Mathematics(4.0), Introduction to Artificial Intelligence(4.0), Computer Graphics(4.0), Operating Systems(4.0), Compiler(4.0), Organization(4.0), Software Engineering

## Papers

- **sciDataQA: Scientific Dataset Recommendation for Question Answering**  
**Jianzhu Yao\***, Zichun Yu\*(equal contribution), Haihong Tang, Jinxiong Xia, Zequn Liu, Houfeng Wang, Sheng Wang  
*Submitted to The 61st Annual Meeting of the Association for Computational Linguistics. **ACL 2023**.*
- **A Benchmark for Understanding and Generating Dialogue between Characters in Stories**  
**Jianzhu Yao**, Ziqi Liu, Jian Guan, Minlie Huang  
*Submitted to The 61st Annual Meeting of the Association for Computational Linguistics. **ACL 2023**.*
- **Crit2SQL: Dataset and Analysis for Criteria2SQL Generation**  
Yuan Zhou, Hanwen Xu, Addie Chambers, Chuwei Xia, Cai Yang, Sihang Zeng, **Jianzhu Yao**, Tianyin Wang, Mingzhe Wu, Yunwei Zhao, William Howard-Snyder, Mingxin Zhang, Zixuan Liu, Yang Lu, Tao Yu, Sheng Wang  
*Submitted to The 61st Annual Meeting of the Association for Computational Linguistics. **ACL 2023**.*
- **EVA2.0: Investigating Open-Domain Chinese Dialogue Systems with Large-Scale Pre-Training**  
Yuxian Gu, Jiaxin Wen, Hao Sun, Yi Song, Pei Ke, Chujie Zheng, Zheng Zhang, **Jianzhu Yao**, Xiaoyan Zhu, Jie Tang, Minlie Huang  
*Machine Intelligence Research 2022*

## Research Experience

### Open-Vocabulary Multi-Hashtag Generation for Social Media Posts

Oct. 2022 – Present

*Advisor: Prof. Bowen Zhou, Zhou Group, Tsinghua University*

*Beijing, China*

- Implemented a sequence-to-sequence model for multi-hashtag generation, designed evaluation metrics with various granularity, and the approach achieved 56.72% entity precision and 34.58% exact-match precision.
- Constructed and cleaned multimodal e-commerce and social media review datasets including text, images, and videos as well as commodities' text and images, including 47K posts and 229K hashtags.
- Participated in a multimodal workshop and gave a presentation on the interdisciplinary of multimodal machine learning and sparsely-activated models like mixture of experts.

### sciDataQA: Scientific Dataset Recommendation for Question Answering

June 2022 – Oct. 2022

*Advisor: Prof. Sheng Wang, Wang Lab, University of Washington, Seattle*

*Beijing, China(Remote)*

- Constructed a large-scale dataset, sciDataQA, including 244K questions and 43K scientific datasets using the Open Pre-trained Transformer (OPT) and a background-enrich generation approach.
- Proposed a novel recursive retrieval approach with the Unified Medical Language System(UMLS) and implemented GCN for retrieval tree embedding to incorporate more domain information into the recommendation process.
- Implemented various applications of the dataset with citation prediction, and improved existing QA systems in the low-resource setting with our dataset.
- Submitted the article to ACL 2023 as the first author.

### A Benchmark for Dialogue Understanding and Generation in Stories

July 2021 – Mar. 2022

*Advisor: Prof. Minlie Huang, Conversational AI Group, Tsinghua University*

*Beijing, China*

- Designed a benchmark for dialogue understanding and generation between characters in stories, which required a higher standard for understanding character relationships and storylines.
- Built a new story dataset with marked dialogue and speakers and tested the performance of existing baselines on the benchmark to investigate the machine's dialogue understanding and generation abilities.
- Proposed learning explicit character representation to guide generation and understanding. The automatic and manual evaluation revealed that our approach outperforms strong baselines by 30% and can generate more coherent and informative dialogue.
- Submitted the article to ACL 2023 as the first author.

## Open-Domain Chinese Dialogue System EVA1.0 and EVA2.0

July 2021 – Mar. 2022

Advisor: Prof. Minlie Huang, Conversational AI Group, Tsinghua University

Beijing, China

- Collaborated with team members to develop the Open-Domain Chinese Dialogue System EVA1.0 and EVA2.0, and was responsible for designing decoding strategies, and constructing dialogue datasets from a massive story corpus.
- Trained and implemented a contradiction detection classifier for Chinese dialogue systems, designed the regeneration pipeline to avoid inconsistent generation, and validated strategy effectiveness using case studies. The classifier could detect 89.07% of the contradictory responses.
- Accepted by Machine Intelligence Research 2022.

## Mixture of Experts Model and Expert Pruning Technique

Mar. 2022 – June 2022

Advisor: Dr. Tao Ge, Natural Language Computing Group, Microsoft Research Asia

Beijing, China

- Researched computing cost trend patterns in all2all communication by conducting profiling experiments using mixture of experts models with different numbers of experts on each device.
- Explored using expert pruning algorithms on the mixture of experts language models to accelerate computation, lower computation latency, and reduce the GPU memory usage by 20% at a given speed with very little performance loss.

## Research on Rotation Invariance of Image Local Features in Object Reconstruction

2021

Advisor: Prof. Shimin Hu, at Graphics & Geometric Computing Group, Tsinghua University

Beijing China

- Researched rotation invariance of image local features in object reconstruction based on deep learning methods and traditional algorithms, which shows existing neural networks suffer from serious data bias and cannot outperform traditional algorithms on some data with rotation invariance.
- Observed the influence of image transformation (like rotation, style transformation, affine transformation, stretching) on feature extraction and matching in SuperPoint, D2-Net, SuperGlue, and SIFT.

## Open-Source Projects

---

- **Ensemble Learning on Spam Classification(Python, SVM, Decision-Tree, Bagging, AdaBoost):** Incorporated SVM and decision-tree with the ensemble learning methods Bagging and AdaBoost for spam classification.
- **KD-Tree Based Stochastic Progressive Photon Mapping Image Rendering Framework(C++):** Implemented the SPPM algorithm, KD-Tree based Bounding Box for intersection acceleration, 3D scene construction, anti-aliasing, and motion blur.
- **Gym Reservation Script(Python, Selenium, PhantomJS, Shell, Wireshark):** Developed a script using Python and PhantomJS to register for a timeslot at the gym with captured cookies by Wireshark.
- **Enterprise Personnel Permissions Management System(Vue, JavaScript, Front-end):** Developed a front-end service for managing personnel permissions in an enterprise(Kuaishou).
- **Education Platform APP Based on Knowledge Graph(Java, SpringBoot, Android Studio, Full-stack):** Developed the front-end and back-end of the educational app IntelEdu based on knowledge graphs with a semantic similarity classifier(BERT).

## Services and Membership

---

- Reviewer: ACL 2023, EMNLP 2022
- Member, Bodybuilding Team of Tsinghua University, 2022 – Present
- Vice President, Winter Swimming Association of Tsinghua University, 2021 – Present
- Member, Tsinghua University Admission Group in Heilongjiang Province, June 2020
- Member, Student Association for Science and Technology, Dept. CST of Tsinghua University, 2020 – Present
- Member, Swimming Team of CST in Tsinghua University, 2019

## Technical Skills

---

**Programming Languages:** Python, LaTeX, Java, C, C++, Shell, HTML/CSS, Assembly(RISC-V, x86)

**Developer Tools:** VS Code, PyCharm, Git, Docker, Linux, Xcode, Unity Hub, Vim, Android Studio, Vivado, Quartus

**Libraries/Frameworks:** PyTorch, Transformers, Fairseq, pytorch-lightning, spaCy, NumPy, Django, SpringBoot, Vue

**TOEFL iBT:** 105 (Reading 28, Listening 28, Speaking 23, Writing 26)

**GRE:** Verbal 155 (66%), Quantitative 170 (96%), Analytical Writing 4.5 (79%)

## Honors & Awards

---

- Sports Excellence Award, Tsinghua University, 2022
- Winter Swimming Star, Tsinghua University, 2022
- First Prize, 37<sup>th</sup> National Physics Competition for College Students in Beijing, China, Beijing Physical Society, 2021
- Academic Excellence Award, Tsinghua University, 2020
- First prize, Engineering Technology Challenge, Tsinghua University, 2019
- First prize, Province Degree in the National Physics Competition (senior), Chinese Physical Society, 2018
- Second prize, Province Degree in the National Mathematical Competition (senior), Chinese Mathematical Society, 2018