

Guanghui Qin

🌐 <https://gqin.me> ✉ gqin2@jhu.edu 🎓 Google Scholar

Education

Johns Hopkins University

Ph.D. in Computer Science (Advisor: Benjamin Van Durme)

Maryland, US

Aug 2019 – Summer 2024 (expected)

Peking University

B.S. in Physics & Computer Science

Beijing, China

Sept 2015 – Jun 2019

Experience

Microsoft Research Lab (MSR)

Research Intern (Mentor: Corby Rosset)

Washington, US

May 2023 – Aug 2023

- Proposed DODO, a *context compression* method for decoder-only LLMs such as LLaMA.
- Applied DODO to downstream tasks, achieving a *20x compression rate* with minimal performance tradeoff on RAG.

Semantic Machines, Microsoft

Research Intern (Mentor: Anthony Platanios)

Remote, US

May 2022 – Aug 2022

- Studied a new research problem for *user action predictions* and built a dataset from GitHub.
- Implemented *graph neural networks* baseline to learn the references between texts and entities.

Center for Language and Speech Processing, Johns Hopkins University

Visiting Researcher (Mentor: Hongyuan Mei and Jason Eisner)

Maryland, US

Jun 2018 – Oct 2018

- Worked on a *particle smoothing* solution for *neural Hawkes process*, a method for temporal event sequence modeling.
- A framework that enables *neural Hawkes process* to interact with the database through *Datalog*.

Microsoft Research-Asia (MSRA)

Research Intern (Mentor: Jin-Ge Yao and Chin-Yew Lin)

Beijing, China

Nov 2017 – Jun 2018

- Proposed a *Semi-HMMs* based model for grounding natural language to structured data.
- Implemented a demo to induce templates from the corpus for *data-to-text generation*.

Awards

- Best Short Paper Awardee in NAACL 2021 Association for Computational Linguistics (ACL), 2021
- Outstanding Reviewer Association for Computational Linguistics (ACL), 2019
- May 4th Scholarship Peking University, 2016 and 2018
- Silver Medalist Chinese Physics Olympiad (CPhO), 2014

Skills

- Programming languages: Python, Rust, JAVA, and C/C++. Other languages: Shell, \LaTeX , SQL.
- Experience in fine-tuning large language models (LLMs), including distributed training, the use of LoRA, and customized transformer architectures.
- Experience in information retrieval, including retrieval-augmented generation (RAG).
- Machine learning tools: PyTorch, Lightning AI, DeepSpeed, FAISS, and PEFT.
- Network: I host a proxy service (WallacePKU) with more than 13k daily active users. I implement front-end, back-end, proxy protocols, and databases.

Academic Service

I serve as a reviewer for the conferences of NeurIPS (2019 and 2020 as secondary; 2021 to 2023), ICLR (2019 and 2020 as secondary; 2021, 2023, and 2024), ICML (2020 and 2021), ACL (2021), EMNLP (2019 to 2022; *outstanding reviewer award* in 2019), NAACL (2024), AACL (2021), and AKBC (2020 as secondary).

Selected Publications

- Dodo: Dynamic Contextual Compression for Decoder-only LMs.
Guanghui Qin, Corby Rosset, Ethan C Chau, Nikhil Rao, Benjamin Van Durme. In *arXiv*. 2024.
- Researchy Questions: A dataset of multi-perspective, compositional questions for LLM web agents.
Corby Rosset, Ho-Lam Chung, **Guanghui Qin**, Ethan C Chau, Zhuo Feng, Ahmed Hassan Awadallah, Jennifer Neville, Nikhil Rao. In *arXiv*. 2024.
- Streaming Sequence Transduction through Dynamic Compression.
Wenting Tan, Yunmo Chen, Tongfei Chen, **Guanghui Qin**, Haoran Xu, Heidi C Zhang, Benjamin Van Durme, Phillip Koehn. In *arXiv*. 2024.
- Nugget: Neural Agglomerative Embeddings of Text.
Guanghui Qin and Benjamin Van Durme. In *Proceedings of the Conference on International Conference on Machine Learning (ICML)*. 2023.
- The NLP Task Effectiveness of Long-Range Transformers.
Guanghui Qin, Yukun Feng, and Benjamin Van Durme. In *European Chapter of the Association for Computational Linguistics (EACL, oral)*. 2023.
- Learning How to Ask: Querying LMs with Mixtures of Soft Prompts.
Guanghui Qin and Jason Eisner. In *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics (NAACL, short)*. 2021. **Best Short Paper Award**
- Everything is all it takes: A multipronged strategy for zero-shot cross-lingual information extraction.
Mahsa Yarmohammadi, Shijie Wu, Marc Marone, Haoran Xu, Seth Ebner, **Guanghui Qin**, Yunmo Chen, J. Guo, Craig Harman, K. Murray, Aaron S. White, Mark Dredze, and Benjamin Van Durme. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP, oral)*. 2021.
- LOME: Large Ontology Multilingual Extraction.
Patrick Xia*, **Guanghui Qin***, Siddharth Vashishtha, Yunmo Chen, Tongfei Chen, Chandler May, Craig Harman, Kyle Rawlins, Aaron Steven White, and Benjamin Van Durme. In *Proceedings of Conference of the European Chapter of the Association for Computational Linguistics (EACL, demo)*. 2021.
- Iterative Paraphrastic Augmentation with Discriminative Span-based Alignment.
Ryan Culkin, J Edward Hu, Elias Stengel-Eskin, **Guanghui Qin**, and Benjamin Van Durme. In *Transactions of the Association for Computational Linguistics (TACL)*, 9:494-509. 2021.
- Neural Datalog Through Time: Informed Temporal Modeling via Logical Specification.
Hongyuan Mei, **Guanghui Qin**, Minjie Xu, and Jason Eisner. In *Proceedings of the Conference on International Conference on Machine Learning (ICML, oral)*. 2020.
- Imputing Missing Events in Continuous-Time Event Streams.
Hongyuan Mei, **Guanghui Qin**, and Jason Eisner. In *Proceedings of the Conference on International Conference on Machine Learning (ICML, oral)*. 2019.
- Learning Latent Semantic Annotations for Grounding Natural Language to Structured Data.
Guanghui Qin, Jin-Ge Yao, Xuening Wang, Jinpeng Wang, and Chin-Yew Lin. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP, oral)*. 2018.

(* indicates equal contribution)

Last updated on Mar 6, 2024.