# Sitao Cheng

sitao.cheng@uwaterloo.ca | https://sitaocheng.github.io/

## RESEARCH INTEREST

I am passionate about Language Agents, Retrieval-augmented Generation (RAG) and Neural-Symbolic Reasoning. I work on knowledge-intensive reasoning, including structured (Knowledge Base, Table), unstructured (Document) and models parametric knowledge. Currently, I focus on the automatic reward modeling.

## EDUCATION

- **University of Waterloo** *09.2025 - Now*
  *Ph.D. in Computer Science - advised by Prof. Victor Zhong* Waterloo, Canada
- **Nanjing University** *09.2021 - 06.2024*
  *M.S. in Computer Science and Technology - advised by Prof. Yuzhong Qu* Nanjing, China
- **University of Electronic Science and Technology of China** *09.2017 - 06.2021*
  *B.E. in Software Engineering - GPA: 3.98 / 4.00 (Top 3)* Chengdu, China

## PUBLICATIONS *\*EQUAL CONTRIBUTION*

**Conference paper.**

[1] **Understanding the Interplay between Parametric and Contextual Knowledge for Large Language Models**. Workshop (Oral) on ACL, 2025. [link]
**Sitao Cheng**, Liangming Pan, Xunjian Yin, Xinyi Wang, William Yang Wang

[2] **Call Me When Necessary: LLMs can Efficiently and Faithfully Reason over Structured Environments**. ACL, 2024. [link]
**Sitao Cheng**, Ziyuan Zhuang, Yong Xu, Fangkai Yang, Chaoyun Zhang, Xiaoting Qin, Xiang Huang, Ling Chen, Qingwei Lin, Dongmei Zhang, Saravan Rajmohan, Qi Zhang

[3] **QueryAgent: a Reliable and Efficient Reasoning Framework with Environmental Feedback-based Self-Correction**. ACL (Oral), 2024. [link]
Xiang Huang*, **Sitao Cheng***, Shanshan Huang, Jiayu Shen, Yong Xu, Chaoyun Zhang, Yuzhong Qu

[4] **Dynamic Evaluation for Oversensitivity in LLMs**. EMNLP, 2025.
Sophia Xiao Pu, **Sitao Cheng**, Xin Eric Wang, William Yang Wang

[5] **MarkQA: a Large Scale KBQA Dataset with Numerical Reasoning**. EMNLP, 2023. [link]
Xiang Huang, **Sitao Cheng**, Yuheng Bao, Shanshan Huang, Yuzhong Qu

[6] **Question Decomposition Tree for Answering Complex Questions over Knowledge Bases**. AAAI (Oral), 2023. [link]
Xiang Huang, **Sitao Cheng**, Yiheng Shu, Yuheng Bao, Yuzhong Qu

[7] **EfficientRAG: Efficient Retriever for Multi-Hop Question Answering**. EMNLP, 2024. [link]
Ziyuan Zhuang*, Zhiyang Zhang*, **Sitao Cheng**, Fangkai Yang, Jia Liu, Shujian Huang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, Qi Zhang

[8] **Disentangling Memory and Reasoning Ability in Large Language Models**. ACL, 2025. [link]
Mingyu Jin, Weidi Luo, **Sitao Cheng**, Xinyi Wang, Wenyue Hua, Ruixiang Tang, William Yang Wang, Yongfeng Zhang

[9] **RuleArena: A Benchmark for Rule-Guided Reasoning with LLMs in Real-World Scenarios**. ACL, 2025. [link]
Ruiwen Zhou, Wenyue Hua, Liangming Pan, **Sitao Cheng**, Xiaobao Wu, En Yu, William Yang Wang

[10] **Thread: A Logic-Based Data Organization Paradigm for How-To Question Answering with Retrieval Augmented Generation**. EMNLP 2025. [link]
Kaikai An, Fangkai Yang, Liqun Li, Junting Lu, **Sitao Cheng**, Shuzheng Si, Lu Wang, Pu Zhao, Lele Cao, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, Qi Zhang, Baobao Chang

**Preprints.**

[1] **LEDOM: An Open and Fundamental Reverse Language Model**. [link]
Xunjian Yin, **Sitao Cheng**, Yuxi Xie, Xinyu Hu, Li Lin, Xinyi Wang, Liangming Pan, William Yang Wang, Xiaojun Wan

## Research Experience

- **University of Waterloo (R2L Lab)** *09.2025 - Now*
  *Advisor: Prof. Victor Zhong. Role: Ph.D. Student* Waterloo, Canada
  - **Topic**: Exploration of automatic environmental feedback for stronger reasoning ability.

- **University of California, Santa Barbara (NLP Group)** *07.2024 - 06.2025*
  *Advisor: Prof. William Wang. Role: Visiting Research Scholar* Santa Barbara, U.S.A
  - **Topic**: Understanding how effective LLMs leverage parametric knowledge when contextual knowledge is given.
    - ∗ Description: We systematically design various relationships between the two knowledge sources: *supportive, complementary, conflicting and irrelevant.* We introduce a new dataset ECHOQA across scientific, factual and commonsense knowledge, to access models ability of echoing their knowledge given contextual information.
    - ∗ Findings: LLMs consistently **suppress their own knowledge** given the context, regardless of models, knowledge types, the relations between two knowledge sources, and various levels of instructions.
    - ∗ Results: One submission on ARR.

- **Microsoft Research Asia (DKI Group)** *10.2023 - 06.2024*
  *Advisor: Yong Xu, Fangkai Yang, Chaoyun Zhang. Role: Research Intern & Mentor of Junior Interns* Beijing, China
  - **Topic 1**: LLMs reasoning framework over structured environments with retrieval-augmented generation (**Readi**) or neural symbolic reasoning (**QueryAgent**).
    - ∗ Description: With large-scaled and heterogeneous structured environments (e.g. Knowledge Graphs, Tables, Databases, etc), how LLMs can reason both efficiently and faithfully? Our intuition is from humans exploration with real-world environments. We adopt LLMs to either directly maintain a reasoning path (Readi), or step-by-step build a query (QueryAgent), both incorporating pertinent information for correction.
    - ∗ Results: Two publications on ACL 2024.
  - **Topic 2**: Efficient iterative retrieval with soley encoder-based models (**EfficientRAG**) and a new data organization paradigm (**THREAD**) for RAG systems.
    - ∗ Description: For better retrieval, it is crucial to model the link between the chunks. We leverage strong understanding ability of LLMs to reason the link between chunks. We design novel retrieval methods for smaller encoder-based models (EfficientRAG) and re-organize the documents (Thread), to model such link.
    - ∗ Results: One publication on EMNLP 2024. One submission on ARR.
  - **Topic 3**: LLM-based Personalized Assistant with "**SurpriseMe**" interaction by Structured Knowledge Graphs.
    - ∗ Description: LLMs not only answer questions with powerful conversational capabilities, but also provide human beings with emotion and interest assistance tailored to their individual experience.
    - ∗ Results: One submission.

- **Nanjing University (Websoft Lab)** *09.2021 - 06.2024*
  *Advisor: Prof. Yuzhong Qu. Role: Student Researcher* Nanjing, China
  - **Topic 1**: Step-by-step query building (**QueryAgent**) with self-correction based on environmental feedback.
    - ∗ Description: In-context learning generates the query on one go, which is unreliable. While current incremental query-building method suffers from hallucination problems, we introduce a functional toolset with environmental feedback and a zero-shot correction method for both reliability and efficiency.
    - ∗ Results: One publication on ACL 2024.
  - **Topic 2**: A KBQA benchmark (**MarkQA**) requiring both multi-hop and numerical reasoning ability.
    - ∗ Description: We propose NR-KBQA to challenge both reasoning ability over knowledge bases. We build a dataset (MarkQA), scaling automatically to 32k from a small number of seeds. We design PyQL query, a function toolset able to seamless SPARQL conversion, as symbolic reasoning steps, alleviating labeling burden.
    - ∗ Results: One publication on EMNLP 2023.
  - **Topic 3**: A question decomposition method (**QDT**) for better multi-hop reasoning over knowledge bases.
    - ∗ Description: We propose a serializable tree-based structure (QDT) to represent complex questions, which can sufficiently split questions with complex structures. We also propose a two-staged generative based method (Clue-Decipher) to ease the uncontrollable nature of generative LMs.
    - ∗ Results: One publication on AAAI 2023.

## Honors and Awards

- **MCM/ICM H Prize, Outstanding Student of Sichuan Province, Outstanding Student Award** –NJU, UESTC

## Skills

- **Professional Skills:** Popular NLP models (LLM applications, Transformers, attention mechanism, etc.), Pytorch, C++, LaTex, Python, SQL
- **Languages:** TOEFL 106, CET-4 CET-6 Excellent
- **Interests:** Body building (over 6x body weight in the Big 3) , Basketball (member of department team), Swim
- **Social Service:** I serve as a personal assistant for a senior impressionist artist in UC Santa Barbara.