

ZEKUN ZHAO

Department of Computer Science and Engineering
University of California, Santa Cruz
zzhao99@ucsc.edu
zekunzhao.github.io

BIO

My research focuses on large language model inference with **knowledge-based reasoning**. I am particularly interested in developing methods that enhance both the understanding and generation of natural language with an emphasis on **efficient and effective performance**. My anticipated graduation date is 03/2027.

EDUCATION

| | |
|---|---------------------|
| UNIVERSITY OF CALIFORNIA, SANTA CRUZ | Sep 2021 - Present |
| <i>Ph.D. in Natural Language Processing</i> | |
| <i>Advisors: Jeffrey Flanigan</i> | |
| UNIVERSITY OF CALIFORNIA, SANTA CRUZ | Sep 2018 – Mar 2021 |
| <i>Master of Science in Computer Science</i> | |
| UNIVERSITY OF CALIFORNIA, BERKELEY | Dec 2017 – Jun 2018 |
| <i>Exchange Student in the Department of Electrical Engineering and Computer Sciences</i> | |
| NANKAI UNIVERSITY, TIANJIN | Sep 2014 – Jun 2018 |
| <i>Bachelor of Engineering in Intelligent Science and Technology</i> | |

PROJECT

| | |
|---|---------------------|
| FORMAL VERIFICATION OF REASONING MODELS | Aug 2025 - Present |
| <ul style="list-style-type: none">Proposed a neuro-symbolic framework to detect hallucinations in step-by-step LLM reasoning via formal verification.Outlined a step-aware pipeline integrating LLM outputs with the Lean 4 theorem prover; each intermediate claim will be formalized and machine-checked.Defined proof-obligation generation, failure-handling, and feedback loops; unverifiable steps are treated as hallucinations.Status: concept/design phase with literature review and initial Lean 4 specification templates completed; implementation pending. | |
| FAST LLM INFERENCE WITH PARALLEL PROMPTING | Sep 2024 - Jun 2025 |
| <ul style="list-style-type: none">Developing a novel parallel inference method for Transformer Large Language Models(LLMs)Improving the LLMs' inference efficiency without compromising generation qualityOptimizing generation latency and throughput with fast parallel generation for document question answer tasksReducing inference time on various popular datasets (SQUAD, QuAC, DROP) by over 70% to the baseline method | |
| IMPLICIT ROLE RECOGNITION IN DOCUMENT | Sep 2023 - Jun 2024 |
| <ul style="list-style-type: none">Designed a novel prompt method with the knowledge graph for document-level implicit role recognitionConstructed the question answer pairs with predicate-argument relations extracted from PropBankImplemented full-document semantic parsing by incorporating concept coreference and implicit role recognition | |

PUBLICATIONS

- Jon Cai, Kristin Wright-Bettner, Zekun Zhao, Shafiuddin Rehan Ahmed, Abijith Trichur Ramachandran, Jeffrey Flanigan, Martha Palmer, and James Martin. 2025. LiDARR: Linking Document AMRs with Referents Resolvers. In Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations), pages 426–435, Vienna, Austria.