

Sihao Chen

Contact: sihaochen96@gmail.com
Website: www.seas.upenn.edu/~sihaoc/

Last Updated: June 4, 2024
GitHub: www.github.com/schen149

OBJECTIVE

Looking for Machine Learning Engineer/Scientist positions specializing in Large Language Models, Retrieval Augmented Generation, and related areas.

EDUCATION

Ph.D. student in Computer and Information Science 06/2019 - now
University of Pennsylvania GPA: 3.98
Thesis Advisor: Dan Roth
Expected Graduation Date: Summer 2024

B.S. in Computer Engineering with Honor 08/2014 - 01/2018
University of Illinois at Urbana Champaign GPA: 3.60

WORK EXPERIENCE

Part-Time Student Researcher 09/2021 - 05/2023, 09/2023 - 01/2024
Google Research
Host: Alex Fabrikant
Topic: Factuality Estimation of Language Model outputs.

Research Intern 05/2023 - 08/2023
Tencent AI Lab
Host: Hongming Zhang, Dong Yu
Topic: Retrieval Augmented Generation [A4, A5], Stability of RLHF Training [P12].

Research Intern 05/2021 - 08/2021, 05/2022 - 08/2022
Google Research
Host: Alex Fabrikant, Donald Metzler
Topic: Data Search Infrastructure + Text Generation [P9, P10]

Research Intern 05/2020 - 08/2020
Google Ads
Host: Kazuo Sone
Topic: Hallucination in Language Models [P7].

PUBLICATIONS

See [Google Scholar](#) page for up-to-date publications + preprints.

Selected Preprints

- [A4] “Beyond Relevance: Evaluate and Improve Retrievers on Perspective Awareness”
Xinran Zhao, Tong Chen, Sihao Chen, Hongming Zhang, Tongshuang Wu.
(In Submission), 2024
- [A3] “Conceptual and Unbiased Reasoning in Language Models”
Ben Zhou, Hongming Zhang, **Sihao Chen**, Dian Yu, Hongwei Wang, Baolin Peng, Dan Roth,
Dong Yu.
(In Submission), 2024
- [A2] “Dense X Retrieval: What Retrieval Granularity Should We Use?”
Tong Chen, Hongwei Wang, **Sihao Chen**, Wenhao Yu, Kaixin Ma, Xinran Zhao, Dong Yu,

Hongming Zhang.
(In Submission), 2023

- [A1] “Towards Corpus-Scale Discovery of Selection Biases in News Coverage: Comparing What Sources Say About Entities as a Start”
Sihao Chen, William Bruno, Dan Roth.
(In Submission), 2023

Peer-Reviewed Publications

- [P15] “The Language Barrier: Dissecting Safety Challenges of LLMs in Multilingual Contexts”
Lingfeng Shen, Weiting Tan, **Sihao Chen**, Yunmo Chen, Jingyu Zhang, Haoran Xu, Boyuan Zheng, Philipp Koehn, Daniel Khashabi
To appear in ACL (Findings), 2024
- [P14] “Sub-Sentence Encoder: Contrastive Learning of Propositional Semantic Representations”
Sihao Chen, Hongming Zhang, Tong Chen, Ben Zhou, Wenhao Yu, Dian Yu, Baolin Peng, Hongwei Wang, Dan Roth, Dong Yu.
To appear in NAACL, 2024
- [P13] “ExpertQA: Expert-Curated Questions and Attributed Answers”
Chaitanya Malaviya, Subin Lee, **Sihao Chen**, Elizabeth Sieber, Mark Yatskar, Dan Roth.
To appear in NAACL, 2024
- [P12] “The Trickle-down Impact of Reward (In-) Consistency on RLHF”
Lingfeng Shen, **Sihao Chen**, Linfeng Song, Lifeng Jin, Baolin Peng, Haitao Mi, Daniel Khashabi, Dong Yu
In ICLR, 2024
- [P11] “Using LLM for Improving Key Event Discovery: Temporal-Guided News Stream Clustering with Event Summaries”
Nishanth Nakshatri, Siyi Liu, **Sihao Chen**, Daniel J. Hopkins, Dan Roth, Dan Goldwasser
In EMNLP Findings, 2023
- [P10] “PropSegEnt: A Large-Scale Corpus for Proposition-Level Segmentation and Entailment Recognition”
Sihao Chen, Senaka Buthpitiya, Alex Fabrikant, Dan Roth, Tal Schuster
In ACL Findings, 2023
- [P9] “Stretching Sentence-pair NLI Models to Reason over Long Documents and Clusters”
Tal Schuster, **Sihao Chen**, Senaka Buthpitiya, Alex Fabrikant, Donald Metzler
In EMNLP Findings, 2022
- [P8] “Design Challenges for a Multi-Perspective Search Engine”
Sihao Chen*, Siyi Liu*, Xander Uyttendaele, Yi Zhang, William Bruno, and Dan Roth
In NAACL Findings, 2022
- [P7] “Improving Faithfulness in Abstractive Summarization with Contrast Candidate Generation and Selection”
Sihao Chen, Fan Zhang, Kazuo Sone and Dan Roth
In NAACL, 2021
- [P6] “MULTIOPED: A Corpus of Multi-Perspective News Editorials”
Siyi Liu, **Sihao Chen**, Xander Uyttendaele and Dan Roth
In NAACL, 2021

- [P5] “Evaluating NLP Models via Contrast Sets”
Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, **Sihao Chen**, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hanna Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang and Ben Zhou.
In EMNLP Findings, 2020
- [P4] “Navigating Information Pollution: Penn’s COVID-19 Information Platform”
Sihao Chen*, Xander Uyttendaele* and Dan Roth.
The Responsible Data Summit, **Spotlight Award**, 2020, [Media Coverage by PennToday](#)
- [P3] “Do VQA Models Know What To Look At?”
Weiyu Du, **Sihao Chen**, Yijie Zhao, Jianbo Shi and Dan Roth.
Women in CV (WiCV) Workshop, ECCV, 2020
- [P2] “PerspectroScope: A Window to the World of Diverse Perspectives”
Sihao Chen, Daniel Khashabi, Chris Callison-Burch and Dan Roth.
ACL (Demo), 2019
- [P1] “See Things from a Different Angle: Discovering Diverse Perspectives about Claims”
Sihao Chen, Daniel Khashabi, Wenpeng Yin, Chris Callison-Burch and Dan Roth.
NAACL, 2019

Book Chapters

- [B1] “Toward Automatic Discovery of Diverse Perspectives”
with Daniel Khashabi and Dan Roth
in “*Creating a More Transparent Internet: The Perspective Web*”
Cambridge University Press, 2022

INVITED TALKS

1. “Language Models as Generative Search Engines – Past, Present and Future.”
At Microsoft Research, 04/2024
2. “Towards a more contextualized view of the web.”
At Google Research, 04/2024
At Allen Institute for Artificial Intelligence, 04/2024
At Duolingo, 03/2024
3. “Towards a more contextualized view of the web.”
At Duolingo, 03/2024
4. “Dense X Retrieval – How We Represent Non-Parametric Knowledge Matters.”
At Google Research, 01/2024
5. “Sub-Sentence Encoder: Contrastive Learning of Propositional Semantic Representation”
At Google Research, Mountain View, CA, 11/2023
6. “Towards Automatic Discovery of Diverse Perspectives”
At *Mid-Atlantic Student Colloquium on Speech, Language and Learning*
University of Maryland College Park, 03/2020

TEACHING

As *Teaching Assistant* –
Spring 2022, *Applied Machine Learning*, Instructor: Dinesh Jayaraman & Mark Yatskar
Spring 2020, *Computational Linguistics*, Instructor: Chris Callison-Burch

PROFESSIONAL SERVICE

As Program Committee Member – AAAI (2020-2022), EMNLP (2019-2023), ACL (2021-2023, 2019), IJCAI (2020), NAACL (2019-2023), Computational Linguistics (2019), ICLR (2024)

NOTABLE OPEN-SOURCE SOFTWARE CONTRIBUTIONS

- “Sub-Sentence Encoder”
General-purpose, multi-vector contextual text embedding model.
<https://github.com/schen149/sub-sentence-encoder>
- “Proposition-based Retrieval”
Smarter dataset indexing for better information retrieval. Featured in open-source libraries for retrieval augmented generation with LLMs, e.g. LangChain and LlamaIndex.
<https://github.com/chentong0/factoid-wiki>