

# Seungju Han

Incoming Ph.D. Student @ Stanford Computer Science  
Research Intern @ NVIDIA

Email: [seungju@stanford.edu](mailto:seungju@stanford.edu)  
Home: <https://seungjuhan.me>

## Education

### Stanford University

Ph.D. Student in Computer Science

Sep, 2025—  
Stanford, CA

### Seoul National University

B.S. in Electrical and Computer Engineering

\* includes a 3-year mandatory military service in South Korea.

Mar, 2017—Aug, 2024\*  
Seoul, Korea

### Seoul Science High School

Specialized high school for students talented in math and science.

Mar, 2014—Feb, 2017  
Seoul, Korea

## Experiences

### NVIDIA

Pre-doctoral Research Intern

Advisor: Prof. Yejin Choi

Nov, 2024—Present  
Santa Clara, CA (Remote)

### Yonsei University, MIR Lab

Visiting Researcher

Advisor: Prof. Youngjae Yu

Mar, 2023—Present  
Seoul, Korea

### Allen Institute for AI, Mosaic Team

Visiting Researcher

Advisor: Prof. Yejin Choi, Dr. Nouha Dziri, Dr. Jack Hessel, Prof. Youngjae Yu

Sep, 2022—Aug, 2024  
Seattle, WA (Remote)

### Hyperconnect

Machine Learning Engineer

Collaborator: Prof. Buru Chang, Dr. Dongyoung Kim, Dr. Sungjoo Ha

Apr, 2019—Aug, 2022  
Seoul, Korea

## Publications

Google Scholar: [https://scholar.google.com/citations?hl=en&user=g\\_anRqAAAAAJ#](https://scholar.google.com/citations?hl=en&user=g_anRqAAAAAJ#)

\* indicates equal contribution.

## Preprints

3. Prismatic Synthesis: Gradient-based Data Diversification Boosts Generalization in LLM Reasoning  
Jaehun Jung, **Seungju Han\***, Ximing Lu\*, Skyler Hallinan\*, Shrimai Prabhumoye, Mostafa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Yejin Choi  
[blog](#)
2. Retro-Search: Exploring Untaken Paths for Deeper and Efficient Reasoning  
Ximing Lu\*, **Seungju Han\***, David Acuna\*, Hyunwoo Kim\*, Jaehun Jung\*, Shrimai Prabhumoye, Niklas Muennighoff, Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Yejin Choi  
[pdf](#)
1. Nemotron-H: A Family of Accurate and Efficient Hybrid Mamba-Transformer Models  
NVIDIA Team, contributed to build pretraining data  
[pdf](#)

## Conference Papers

16. MAPoRL: Multi-Agent Post-Co-Training for Collaborative Large Language Models with Reinforcement Learning  
Chanwoo Park, **Seungju Han**, Xingzhi Guo, Asuman Ozdaglar, Kaiqing Zhang, Joo-Kyung Kim  
**ACL 2025**, [pdf](#)
15. Representation Bending for Large Language Model Safety  
Ashkan Yousefpour\*, Taeheon Kim\*, Ryan S Kwon, Seungbeen Lee, Wonje Jeung, **Seungju Han**, Alvin Wan, Harrison Ngan, Youngjae Yu, Jonghyun Choi  
**ACL 2025**, [pdf](#)
14. AI as Humanity’s Salieri: Quantifying Linguistic Creativity of Language Models via Systematic Attribution of Machine Text against Web Text  
Ximing Lu, Melanie Sclar, Skyler Hallinan, Niloofar Mireshghallah, Jiacheng Liu, **Seungju Han**, Allyson Ettinger, Liwei Jiang, Khyathi Chandu, Nouha Dziri, Yejin Choi  
**ICLR 2025 (Oral Presentation)**, [pdf](#)
13. Do LLMs Have Distinct and Consistent Personality? TRAIT: Personality Testset designed for LLMs with Psychometrics  
Seungbeen Lee\*, Seungwon Lim\*, **Seungju Han**, Giyoung Oh, Minju Kim, Beongwoo Kwak, Jiwan Chung, Hyunjoo Chae, Dongha Lee, Jinyoung Yeo, Youngjae Yu  
**Findings of NAACL (long) 2025**, [code](#), [pdf](#)
12. WILDGUARD: Open One-stop Moderation Tools for Safety Risks, Jailbreaks, and Refusals of LLMs  
**Seungju Han\***, Kavel Rao\*, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, Nouha Dziri  
**NeurIPS 2024 Datasets & Benchmarks**, [code](#), [pdf](#)
11. WILDTEAMING at Scale: From In-the-Wild Jailbreaks to (Adversarially) Safer Language Models  
Liwei Jiang, Kavel Rao\*, **Seungju Han\***, Faeze Brahman, Sachin Kumar, Niloofar Mireshghallah, Ximing Lu, Marteen Sap, Yejin Choi, Nouha Dziri  
**NeurIPS 2024**, [code](#), [pdf](#)
10. Selective Vision is the Challenge for Visual Reasoning: A Benchmark for Visual Argument Understanding  
Jiwan Chung, Sungje Lee, Minseo Kim, **Seungju Han**, Ashkan Yousefpour, Jack Hessel, Youngjae Yu  
**EMNLP 2024 (long) (Oral Presentation)** [code](#), [pdf](#)
9. Multimodal Laughter Reasoning with Textual Audio-Visual Representation  
Hyun Lee, Sung Bin Kim, **Seungju Han**, Youngjae Yu, Tae Hyun Oh  
**Findings of NAACL (long) 2024**, [code](#), [pdf](#)  
**ICCV Workshop, What is Next in Multimodal Foundation Models? 2023**
8. Reading Books is Great, But Not if You Are Driving! Visually Grounded Reasoning about Defeasible Commonsense Norms  
**Seungju Han**, Junhyeok Kim, Jack Hessel, Liwei Jiang, Jiwan Chung, Yejin Son, Yejin Choi, Youngjae Yu  
**EMNLP 2023 (long) (Oral Presentation)**, [code](#), [project page](#), [pdf](#)
7. CHAMPAGNE: Learning Real-world Conversation from Large-Scale Web Videos  
**Seungju Han**, Jack Hessel, Nouha Dziri, Yejin Choi, Youngjae Yu  
**ICCV 2023**, [code](#), [project page](#), [pdf](#)
6. Measuring and Improving Semantic Diversity of Dialogue Generation  
**Seungju Han**, Beomsu Kim, Buru Chang  
**Findings of EMNLP (long) 2022**, [code](#), [pdf](#)
5. Meet Your Favorite Character: Open-domain Chatbot Mimicking Fictional Characters with only a Few Utterances  
**Seungju Han\***, Beomsu Kim\*, Jin Yong Yoo\*, Seokjun seo, Sangbum Kim, Enkhbayar Erdenee, Buru Chang  
**NAACL (short) 2022**, [code](#), [pdf](#)

4. Understanding and Improving the Exemplar-based Generation for Open-domain Conversation  
**Seungju Han\***, Beomsu Kim\*, Seokjun seo\*, Enkhbayar Erdenee\*, Buru Chang  
**ACL 4th Workshop on NLP4ConvAI (Oral Presentation, Outstanding Paper) 2022**, [code](#), [pdf](#)
3. Distilling the Knowledge of Large-scale Generative Models into Retrieval Models for Efficient Open-domain Conversation  
 Beomsu Kim\*, Seokjun seo\*, **Seungju Han\***, Enkhbayar Erdenee\*, Buru Chang  
**Findings of EMNLP (long) 2021**, [code](#), [pdf](#)
2. Disentangling Label Distribution for Long-tailed Visual Recognition  
 Youngkyu Hong\*, **Seungju Han\***, Kwanghee Choi\*, Seokjun seo, Beomsu Kim, Buru Chang  
**CVPR 2021**, [code](#), [blog](#), [pdf](#)
1. Attention: Few-Shot Text-to-Speech Utilizing Attention-Based Variable-Length Embedding  
 Seungwoo Choi\*, **Seungju Han\***, Dongyoung Kim\*, Sungjoo Ha  
**Interspeech 2020**, [pdf](#)

## Awards and Honors

Ph.D. Scholarship from <b>Korea Foundation for Advanced Studies (KFAS)</b> Doctoral study fellowship, declined since started Ph.D. a year later	2023—2024
Undergraduate Scholarship from <b>Kwanjeong Educational Foundation</b> Full scholarship (awarded to up to 50 people every year)	2019—2023
Undergraduate Scholarship from <b>Seoul National University</b> Full scholarship (merit-based)	2017—2019
Scholarship from <b>Hanseong Son Jae Han Scholarship Foundation</b> Hanseong Nobel scholarship (awarded to up to 200 people every year)	2015—2016

## Professional Activities

### Volunteer

EMNLP	2022
-------	------

### Reviewer (\* indicates outstanding reviewer)

EMNLP	2021
ACL ARR	2022, 2023, 2024, 2025
NeurIPS	2023, 2024*
ICLR	2024
ICML	2025