

Seungju Han

Email: wade3han@snu.ac.kr Homepage: <https://seungjuhan.me> Twitter/X: @SeungjuHan3

Education

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

University of Washington

Exchange Student

Sep, 2022—Dec, 2022
Seattle, WA

Seoul Science High School

Specialized high school for students talented in math and science.

Mar, 2014—Feb, 2017
Seoul, Korea

Experiences

University of Washington, xlab

Undergraduate Research Intern

Advisor: Prof. Yejin Choi

Aug, 2024—Present
Seattle, WA (Remote)

Yonsei University, MIR Lab

Visiting Researcher

Advisor: Prof. Youngjae Yu

Mar, 2023—Present
Seoul, Korea

Allen Institute for AI, Mosaic Team

Visiting Predoctoral Researcher

Advisor: Prof. Yejin Choi / Hosts: 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 and Preprints

Google Scholar: https://scholar.google.com/citations?hl=en&user=g_anRqAAAAAJ#

* indicates equal contribution.

Preprints

2. 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
Under review
1. 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, Hyungjoo Chae, Dongha Lee, Jinyoung Yeo, Youngjae Yu
Preprint, [code](#), [pdf](#)

Publications

12. WILDDGUARD: 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) [code](#), [pdf](#)
9. 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](#)
8. 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
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. ATTENTRON: 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

Research Grant from Korea Foundation for Advanced Studies (KFAS) Doctoral study fellowship (awarded to six people this year in Computer Science)	2023—2024
Undergraduate Scholarship from Kwanjeong Educational Foundation Full scholarship, approx. 20,000 USD in total (awarded to up to 50 people every year)	2019—2023
Undergraduate Scholarship from Seoul National University Full scholarship, approx. 9,000 USD in total (merit-based)	2017—2019

Scholarship from **Hanseong Son Jae Han Scholarship Foundation**
Hanseong Nobel scholarship, approx. 10,000 USD in total (awarded to up to 200 people every year)

2015—2016

Invited Presentations

Learning representations from large-scale videos, Tech seminar @ Dalphi AI
A path towards AGI, Tech seminar @ Corca AI
Research in start-up, Tech seminar @ FriendliAI
Research in start-up, Tech seminar @ Corca AI
Open-domain conversation agent, AI retreat @ Seoul National University

2023 Summer
2023 Summer
2022 Fall
2022 Summer
2021 Spring

Professional Activities

Volunteer

EMNLP (Virtual)

2022

Reviewer

EMNLP
ACL ARR
NeurIPS
ICLR

2021
2022—
2023—
2024—